

**Western Alaska Salmon Stock Identification Project
Technical Document 7: Regional Allocation Model**

by

James R. Jasper,

Christopher Habicht,

Andrew R. Munro,

and

William D. Templin

May 2012

Alaska Department of Fish and Game

Divisions of Sport Fish and Commercial Fisheries



Symbols and Abbreviations

The following symbols and abbreviations, and others approved for the *Système International d'Unités* (SI), are used without definition in the following reports by the Divisions of Sport Fish and of Commercial Fisheries: Fishery Manuscripts, Fishery Data Series Reports, Fishery Management Reports, Special Publications and the Division of Commercial Fisheries Regional Reports. All others, including deviations from definitions listed below, are noted in the text at first mention, as well as in the titles or footnotes of tables, and in figure or figure captions.

Weights and measures (metric)		General		Mathematics, statistics	
centimeter	cm	Alaska Administrative Code	AAC	<i>all standard mathematical signs, symbols and abbreviations</i>	
deciliter	dL	all commonly accepted abbreviations	e.g., Mr., Mrs., AM, PM, etc.	alternate hypothesis	H_A
gram	g	all commonly accepted professional titles	e.g., Dr., Ph.D., R.N., etc.	base of natural logarithm	e
hectare	ha	at	@	catch per unit effort	CPUE
kilogram	kg	compass directions:		coefficient of variation	CV
kilometer	km	east	E	common test statistics	(F, t, χ^2 , etc.)
liter	L	north	N	confidence interval	CI
meter	m	south	S	correlation coefficient	
milliliter	mL	west	W	(multiple)	R
millimeter	mm	copyright	©	correlation coefficient (simple)	r
		corporate suffixes:		covariance	cov
Weights and measures (English)		Company	Co.	degree (angular)	$^\circ$
cubic feet per second	ft ³ /s	Corporation	Corp.	degrees of freedom	df
foot	ft	Incorporated	Inc.	expected value	E
gallon	gal	Limited	Ltd.	greater than	>
inch	in	District of Columbia	D.C.	greater than or equal to	\geq
mile	mi	et alii (and others)	et al.	harvest per unit effort	HPUE
nautical mile	nmi	et cetera (and so forth)	etc.	less than	<
ounce	oz	exempli gratia	e.g.	less than or equal to	\leq
pound	lb	(for example)		logarithm (natural)	ln
quart	qt	Federal Information Code	FIC	logarithm (base 10)	log
yard	yd	id est (that is)	i.e.	logarithm (specify base)	\log_2 , etc.
		latitude or longitude	lat. or long.	minute (angular)	'
Time and temperature		monetary symbols (U.S.)	\$, ¢	not significant	NS
day	d	months (tables and figures): first three letters	Jan, ..., Dec	null hypothesis	H_0
degrees Celsius	$^\circ\text{C}$	registered trademark	®	percent	%
degrees Fahrenheit	$^\circ\text{F}$	trademark	™	probability	P
degrees kelvin	K	United States (adjective)	U.S.	probability of a type I error (rejection of the null hypothesis when true)	α
hour	h	United States of America (noun)	USA	probability of a type II error (acceptance of the null hypothesis when false)	β
minute	min	U.S.C.	United States Code	second (angular)	"
second	s	U.S. state	use two-letter abbreviations (e.g., AK, WA)	standard deviation	SD
Physics and chemistry				standard error	SE
all atomic symbols				variance	
alternating current	AC			population sample	Var
ampere	A			sample	var
calorie	cal				
direct current	DC				
hertz	Hz				
horsepower	hp				
hydrogen ion activity (negative log of)	pH				
parts per million	ppm				
parts per thousand	ppt, ‰				
volts	V				
watts	W				

REGIONAL INFORMATION REPORT 5J12-12

**WESTERN ALASKA SALMON STOCK IDENTIFICATION PROJECT
TECHNICAL DOCUMENT 7: REGIONAL ALLOCATION MODEL**

by

James R. Jasper, Christopher Habicht, Andrew R. Munro and William D. Templin
Alaska Department of Fish and Game, Division of Commercial Fisheries, Gene Conservation Laboratory,
Anchorage

Alaska Department of Fish and Game
Division of Sport Fish, Research and Technical Services
333 Raspberry Road, Anchorage, Alaska, 99518-1565

May 2012

The Regional Information Report Series was established in 1987 and was redefined in 2006 to meet the Division of Commercial Fisheries regional need for publishing and archiving information such as project operational plans, area management plans, budgetary information, staff comments and opinions to Board of Fisheries proposals, interim or preliminary data and grant agency reports, special meeting or minor workshop results and other regional information not generally reported elsewhere. Reports in this series may contain raw data and preliminary results. Reports in this series receive varying degrees of regional, biometric and editorial review; information in this series may be subsequently finalized and published in a different department reporting series or in the formal literature. Please contact the author or the Division of Commercial Fisheries if in doubt of the level of review or preliminary nature of the data reported. Regional Information Reports are available through the Alaska State Library and on the Internet at <http://www.adfg.alaska.gov/sf/publications/>.

Note: This document serves as a record of communication between the Alaska Department of Fish and Game Commercial Fisheries Division and the Western Alaska Salmon Stock Identification Program Technical Committee. As such, these documents serve diverse ad hoc information purposes and may contain basic, uninterpreted data. The contents of this document have not been subjected to review and should not be cited or distributed without the permission of the authors or the Commercial Fisheries Division.

Note: The appearance of product names or specific company names is not an Alaska Department of Fish and Game recommendation for or implied endorsement. The Alaska Department of Fish and Game, in accordance with State of Alaska ethics laws, does not favor one group over another through endorsement or recommendation.

*James R. Jasper, Christopher Habicht, Andrew R. Munro, and William D. Templin
Alaska Department of Fish and Game, Division of Commercial Fisheries,
333 Raspberry Road, Anchorage, AK, 99518-1565 USA*

This document should be cited as:

Jasper, J. R., C. Habicht, A. R. Munro, and W. D. Templin. 2012. Western Alaska Salmon Stock Identification Project Technical Document 7: Regional allocation model. Alaska Department of Fish and Game, Division of Commercial Fisheries, Regional Information Report 5J12-12, Anchorage.

The Alaska Department of Fish and Game (ADF&G) administers all programs and activities free from discrimination based on race, color, national origin, age, sex, religion, marital status, pregnancy, parenthood, or disability. The department administers all programs and activities in compliance with Title VI of the Civil Rights Act of 1964, Section 504 of the Rehabilitation Act of 1973, Title II of the Americans with Disabilities Act (ADA) of 1990, the Age Discrimination Act of 1975, and Title IX of the Education Amendments of 1972.

If you believe you have been discriminated against in any program, activity, or facility please write:

ADF&G ADA Coordinator, P.O. Box 115526, Juneau, AK 99811-5526

U.S. Fish and Wildlife Service, 4401 N. Fairfax Drive, MS 2042, Arlington, VA 22203

Office of Equal Opportunity, U.S. Department of the Interior, 1849 C Street NW MS 5230, Washington DC 20240

The department's ADA Coordinator can be reached via phone at the following numbers:

(VOICE) 907-465-6077, (Statewide Telecommunication Device for the Deaf) 1-800-478-3648,

(Juneau TDD) 907-465-3646, or (FAX) 907-465-6078

For information on alternative formats and questions on this publication, please contact:

ADF&G, Division of Sport Fish, Research and Technical Services, 333 Raspberry Rd, Anchorage AK 99518 (907) 267-2375

TABLE OF CONTENTS

	Page
LIST OF TABLES.....	i
LIST OF FIGURES.....	ii
ABSTRACT.....	1
INTRODUCTION.....	1
METHODS.....	2
General Bayesian Methods.....	3
The Pella-Masuda Model.....	3
Regional Allocation Model.....	5
Simulations.....	8
RESULTS.....	8
DISCUSSION.....	9
ACKNOWLEDGEMENTS.....	10
REFERENCES CITED.....	11
TECHNICAL COMMITTEE REVIEW AND COMMENTS.....	12
ADDITIONAL COMMENTARY.....	15
FIGURES.....	17

LIST OF TABLES

Table	Page
1. Simulation results and root mean square error (rMSE) for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model with the True Flat Prior (P-M TFP), the Pella-Masuda Model with the Regional Flat Prior (P-M RFP), and the Regional Allocation Model (RAM).....	8

LIST OF FIGURES

Figure	Page
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. Regional means are indicated.	18
2. Unweighted pair-group method (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.	19
3. Simulation results for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model using the True Flat Prior. The height of the bars represents the mean of 100 repetitions with the 90% credibility interval indicated. The horizontal rule is 90% correct allocation. Numbers under labels are the number of stocks within the region. These results are the same as shown in Figure 1 with the stock proportions summed into regions.	20
4. Simulation results for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model using the True Flat Prior. The height of the bars represents the mean of 100 repetitions with the 90% credibility interval indicated. The horizontal rule is 90% correct allocation. Numbers under labels are the number of stocks within the region. Numbers under labels are the number of stocks within the region.....	21
5. Simulation results for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model using the True Flat Prior. The height of the bars represents the mean of 100 repetitions with the 90% credibility interval indicated. The horizontal rule is 90% correct allocation. Numbers under labels are the number of stocks within the region.	22

ABSTRACT

Uncertainty about the magnitude, frequency, location, and timing of the nonlocal harvest of sockeye and chum salmon in Western Alaska fisheries was the impetus for the Western Alaska Salmon Stock Identification Project (WASSIP). The project was designed to use genetic data in mixed stock analysis to reduce this uncertainty. Mixed stock analysis methods for estimating stock (population) compositions in fisheries have evolved over time, but the recent “gold-standard” has been the Pella-Masuda Bayesian model. Recent observations in the Gene Conservation Laboratory indicate that disproportionate numbers of stocks within a region can lead to significant bias in regional composition estimates when regional stock structure is shallow. We present a new analytical model that appears to diminish this bias. Baseline data for chum salmon stocks from Western Alaska were selected because these stocks represent weak regional structure. 100 mixtures of 200 fish, each composed of 100% Norton Sound chum were analyzed with a Western Alaska baseline using three methods: 1) the Pella-Masuda Model with True Flat Prior (TFP); 2) Pella-Masuda Model using Regional Flat Prior (RFP), and the third method, termed the Regional Allocation Model (RAM), currently under development at the Gene Conservation Laboratory. Posterior means of the stock proportions and the regional proportions were calculated as well as the means, central 90% quantiles, and root mean square errors of the 100 posterior means. The mean and central 90% of the Norton Sound proportions were 0.831 (0.686–0.929) for the TFP, 0.834 (0.696–0.932) for the RFP, and 0.880 (0.757–0.949) for the RAM; and the root mean square errors were 0.091 for the TFP, 0.088 for the RFP, and 0.063 for the RAM. For the Pella-Masuda Model, both the TFP and the RFP showed very similar amounts of misallocation, but the RFP tended to shift some of the misallocation away from the regions with the most stocks and into regions with fewer stocks. The RAM showed less misallocation than both prior structures of the Pella-Masuda Model in terms of point estimate and tightness of the central 90% quantiles, and tended to flatten out the amount of misallocation more evenly across the remaining regions. The RAM appeared to be moderately successful in reducing the nonuniform bias due to the unequal distribution in the number of stocks among the regions. The RAM presented here is extended to only two levels of hierarchy of stocks within regions; it is conceivable to expand this model to further levels of hierarchy, such as substocks within stocks, and stocks within regions to be useful in situations where multiple levels of structure exist.

Key words: Western Alaska Salmon Stock Identification Project, WASSIP, mixed stock analysis, MSA, Pella-Masuda Model, Bayesian analytical methods, Regional Allocation Model, chum salmon

INTRODUCTION

Mixed stock analysis methods for estimating stock (population) compositions in fisheries have evolved over time from conditional maximum-likelihood (Fournier et al. 1984) to Bayesian (Pella and Masuda, 2001) approaches. The Pella-Masuda Model (a Bayesian approach; Pella and Masuda, 2001) has been the “gold standard” since 2001. In these methods, however, bias is inevitable because the estimation of the stock proportions is constrained to be nonnegative and sum to one, meaning that rare or absent stocks in the mixture are overestimated while common stocks are under estimated (Pella and Milner, 1987). Stocks are usually grouped into regional stock groupings (regions) for reporting.

Recent observations in our laboratory indicate that disproportionate numbers of stocks within a region can lead to significant bias in regional composition estimates when regional stock structure is shallow. We have observed that regions represented by large numbers of stocks seem to acquire higher misallocations than regions represented by fewer stocks (Figure 1). This bias can be reduced at the regional level by grouping stocks with similar genetic attributes into regions, then summing estimated proportions across stocks within the regions (Wood et al. 1987). Here we present a rationale for why we think the observed nonuniform bias occurred and a method that appears to improve allocation at the regional level as well as distribute the misallocation more evenly among regions.

In the Pella-Masuda Model, the data augmentation algorithm is used to generate from the posterior distribution the stock identities of each of the mixture individuals, and then generate the

stock proportions and baseline allele frequencies based on summaries of these identities. At each cycle of the algorithm, the stock identity of mixture individual m is stochastically assigned to stock i with probability proportional to the product of stock i 's contribution to the mixture and the relative frequency of individual m 's genotype in stock i . This means that individual m has a finite probability of belonging to each and every stock in the baseline.^a We will refer to these probabilities as the identity probabilities.

The chances that individual m is assigned to the correct stock at a particular iteration is a function of not only the genetic distinction of its stock, but also, heuristically, the number of stocks in the baseline. Fortunately, fisheries managers often are not interested in the proportion of individual stocks, but rather in the contribution made by all stocks within regions. If the stocks within a region are genetically more similar to each other than to stocks in other regions (strong regional structure), then the chances of correctly assigning an individual to a stock within the correct region each cycle greatly improves estimation (Wood et al. 1987). However, with weak regional structure, the chances of assigning an individual to a stock within the correct region may be significantly influenced by the number of stocks in each region. This may be because the probability of assigning an individual to a particular region is the sum of the identity probabilities across all the stocks in the region, such that adding stocks adds probability. If the amount of misallocation to a region is a function of the number of stocks within that region, an inherent nonuniform bias in regional contribution estimates can occur simply due to differing numbers of stocks among regions.

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

METHODS

We considered 3 methods to examine the assertion that unequal numbers of stocks within regions do not affect bias in misallocation.^b We selected baseline data for chum salmon stocks from Western Alaska. These data were chosen because these stocks represent weak regional structure (Figure 2).

The first 2 methods use the Pella-Masuda Model but differ in how the priors are assigned. The first method is the widely used True Flat Prior (TFP; Pella and Masuda, 2001). This model provides no *a priori* information about the regional structure and gives an equal prior “count” of $1/C$ to each of the stocks in the baseline, where C is the number of stocks. This is the model that provided the recent observations in our laboratory that suggested that disproportionate numbers of stocks within a reporting group can affect the regional composition estimates.

The second method, termed the Regional Flat Prior (RFP), is a method currently in use at ADF&G's Gene Conservation Laboratory (Dann et al. 2009). The structure of the prior for stock proportions is an *ad hoc* alternative to the TFP. Under the RFP, for each of the stocks within the g th region, we give a prior “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 region. Therefore, equal prior “count” is given to each

^a This sentence is commented on in the section entitled “Technical Committee Review and Comments.”

^b This sentence is commented on in the section entitled “Technical Committee Review and Comments.”

region, but the prior “count” given to a stock is dependent upon the number of stocks within its region.

The third method, termed the Regional Allocation Model (RAM), is currently under development at the Gene Conservation Laboratory. This model is very similar to the Pella-Masuda Model in that it is based on the data augmentation algorithm that alternates between generating the parameters of the model. The difference is that in the RAM, we first generate the regional identity of each individual, and then produce regional contributions based on summaries of these regional identities. For individual m , the regional identity probability of belonging to region g is proportional to region g 's contribution to the mixture times a weighted average relative frequency of individual m 's genotype across all C_g stocks within the region. The weights are simply the within-region stock proportions, and they sum to 1. Because the weights do sum to 1, the genetic component of the regional identity probabilities remain on the same scale regardless of the number of stocks within the region, which should presumably moderate the nonuniform bias due to the unequal distribution of stocks among the regions. There is actually a second stage to the data augmentation algorithm in which, after an individual is assigned to a region, it is then allocated to a stock within that region. This is done exactly as is done in the Pella-Masuda Model except that it is done with respect to a baseline that is reduced to only that region.

General Bayesian Methods

For estimating parameters θ from data X using Bayesian methods, we aim at the evaluation of the posterior distribution $P(\theta|X) = L(X|\theta) P(\theta)/m(X)$, where $L(X|\theta)$ is the likelihood of the data given the parameters, $P(\theta)$ is the prior distribution of the parameters, which must be specified, and $m(X)$ is the constant marginal distribution of the data. From this distribution, summary statistics for θ can be derived. However, these distributions are rarely soluble in closed form for multidimensional parameter vector θ , and we must rely on drawing samples from it via a Gibbs sampling routine, from which the summary statistics can be calculated. For mixed stock analysis, θ represents the stock proportions and the baseline allele frequencies while X corresponds to the mixture genotypes and the baseline allele counts. As mentioned previously, a prior distribution must be specified for the parameters. In the forthcoming models, the mathematically convenient Dirichlet distribution is used for the stock proportions as well as the baseline allele frequencies. A Dirichlet distribution with parameter vector λ is a distribution on a vector W whose sum is constrained to 1. It has the form:

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

The Pella-Masuda Model

We denote the count of the j th ($j=1,2,\dots,J_d$) allele of the d th ($d=1,2,\dots,D$) locus for mixture individual m as x_{mdj} , and let X_m signify the entire multi-locus genotype for this individual. The array X represents the multi-locus genotypes for all M mixture individuals. Similarly, we let y_{idj} denote the count of the j th allele for the d th locus of the i th baseline stock, and Y denotes the entire baseline. This describes the data.

To describe the parameters, let the stock proportion for the i th stock be denoted as P_i , and let \mathbf{P} be the vector of all stock proportions. We place a Dirichlet prior distribution on the stock proportions with prior parameters $\boldsymbol{\alpha}$, where α_i is determined by our choice of prior structure discussed earlier (RFP or TFP).

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 \mathbf{Q} denote the entire array of baseline relative frequencies. We place a Dirichlet prior distribution on \mathbf{Q}_{id} with prior parameters $\boldsymbol{\beta}_d$, where $\beta_{dj} = 1/J_d$, with J_d being the number of alleles for locus d (Rannala-Mountain 1997).

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 one if individual m belongs to the i th stock and zero otherwise. We denote \mathbf{Z}_m as the vector of stock identities for individual m , and \mathbf{Z} as the matrix of stock identities for the entire mixture. We place a multinomial prior on \mathbf{Z}_m with size 1 and probabilities equal to the stock proportions \mathbf{P} .

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

$$f(\mathbf{X}_m|\mathbf{Q}_i) \propto \prod_{d=1}^D \prod_{j=1}^{J_d} q_{idj}^{x_{dj}}$$

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

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

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

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

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

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

$$\begin{aligned} P(\mathbf{P}, \mathbf{Q}, \mathbf{Z}|\mathbf{X}, \mathbf{Y}) &\propto L(\mathbf{X}, \mathbf{Y}|\mathbf{Q}, \mathbf{Z})P(\mathbf{Z}|\mathbf{P})P(\mathbf{P}|\boldsymbol{\alpha})P(\mathbf{Q}|\boldsymbol{\beta}) \\ &\propto \left(\prod_{m=1}^M \prod_{i=1}^C f(\mathbf{X}_m|\mathbf{Q}_i)^{z_{mi}} \right) \left(\prod_{i=1}^C \prod_{d=1}^D \prod_{j=1}^{J_d} q_{idj}^{y_{idj}} \right) \\ &\times \left(\prod_{m=1}^M \prod_{i=1}^C P_i^{z_{mi}} \right) \left(\prod_{i=1}^C P_i^{\alpha_i} \right) \left(\prod_{i=1}^C \prod_{d=1}^D \prod_{j=1}^{J_d} q_{idj}^{\beta_{dj}} \right) \end{aligned}$$

The benefit of using the chosen prior distributions is that the conditional posterior distribution for each of the parameters given the data and the remaining parameters is of the same form as the

prior distribution (conjugacy). This property makes them easy to sample from within a Gibbs sampler, which proceeds as follows: first, starting with initial values for \mathbf{P} and \mathbf{Q} , we draw stock identities for each of the mixture individuals from:

$$\mathbf{Z}_m | \mathbf{P}, \mathbf{Q}, \mathbf{X}_m \sim \text{multinomial} \left(1, \left\{ \frac{P_i f(\mathbf{X}_m | \mathbf{Q}_i)}{\sum_{k=1}^C P_k f(\mathbf{X}_m | \mathbf{Q}_k)} \right\}_{i=1,2,\dots,C} \right)$$

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

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

Finally, for each stock and for each locus, we generate \mathbf{Q}_{id} from:

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

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

Regional Allocation Model

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

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

Denote S_{gk} as the within-region stock proportion for the k th stock in the g th region, and denote \mathbf{S}_g as the vector of all C_g stock proportions within the g th region—again, notice that \mathbf{S}_g must sum to 1. We place a Dirichlet prior distribution on \mathbf{S}_g with parameters $\boldsymbol{\delta}_g$, with $\delta_{gk} = 1/C_g$. The ragged array of all stock proportions is represented by \mathbf{S} .^c

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

We let r_{mg} denote the regional identity for the g th stock for the m th mixture individual, where $r_{mg} = 1$ if individual m belongs to the g th region, and zero otherwise. The vector of regional identities for the m th individual is denoted as \mathbf{r}_m , and the matrix of all regional identities is

^c This sentence is commented on in the section entitled “Technical Committee Review and Comments.”

represented as \mathbf{r} . A multinomial prior distribution is placed on \mathbf{r}_m with size 1 and probabilities equal to the regional contributions \mathbf{R} .

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 individual, where $z_{mgk} = 1$ if individual m belongs to the k th stock of the g th region, and zero otherwise. Denote \mathbf{z}_{mg} as the vector of stock identities for the g th region for the m th individual, and let \mathbf{z}_m be the ragged array of stock identities for this individual. The ragged array of all stock identities is denoted as \mathbf{z} . We place a multinomial prior distribution on \mathbf{z}_{mg} with size r_{mg} and probabilities equal to \mathbf{S}_g . Because $r_{mg} = 1$ if individual m belongs to the g th region, and zero otherwise, the only way the prior distribution of \mathbf{z}_{mg} can have positive size is if $r_{mg} = 1$. In other words, the m th individual cannot belong to a stock that is outside that individual's region.

If we knew both the region and stock of origin for each mixture individual, the full genotypic likelihood can be expressed as:

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

Here, we use $I()$ as an indicator function that is equal to 1 if the argument is true, and zero otherwise. Similar to the previous model, the baseline likelihood can be written as:

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

The full likelihood, $L(\mathbf{X}, \mathbf{Y}|\mathbf{Q}, \mathbf{r}, \mathbf{z})$, is simply the product of these two components. Multiplying the likelihood by the priors gives the posterior distribution:

$$\begin{aligned} P(\mathbf{R}, \mathbf{S}, \mathbf{Q}, \mathbf{r}, \mathbf{z} | \mathbf{X}, \mathbf{Y}) &\propto L(\mathbf{X}, \mathbf{Y} | \mathbf{Q}, \mathbf{r}, \mathbf{z}) P(\mathbf{z} | \mathbf{r}, \mathbf{S}) P(\mathbf{r} | \mathbf{R}) P(\mathbf{S} | \boldsymbol{\delta}) P(\mathbf{R} | \boldsymbol{\gamma}) P(\mathbf{Q} | \boldsymbol{\beta}) \\ &\propto \left\{ \prod_{m=1}^M \prod_{g=1}^G \left(\prod_{k=1}^{C_g} f(\mathbf{X}_m | \mathbf{Q}_{gk})^{z_{mgk}} I \left(\sum_{k=1}^{C_g} z_{mgk} = r_{mg} \right) \right) \right\} \left\{ \prod_{g=1}^G \prod_{k=1}^{C_g} \prod_{d=1}^D \prod_{j=1}^{J_d} q_{gk d j}^{y_{gk d j}} \right\} \\ &\times \left\{ \prod_{m=1}^M \prod_{g=1}^G \left(\prod_{k=1}^{C_g} S_{gk}^{z_{mgk}} I \left(\sum_{k=1}^{C_g} z_{mgk} = r_{mg} \right) \right) \right\} \left\{ \prod_{m=1}^M \prod_{g=1}^G R_g^{r_{mg}} \right\} \\ &\times \left\{ \prod_{g=1}^G \prod_{k=1}^{C_g} S_{gk}^{\delta_{gk}} \right\} \left\{ \prod_{g=1}^G R_g^{\gamma_g} \right\} \left\{ \prod_{g=1}^G \prod_{k=1}^{C_g} \prod_{d=1}^D \prod_{j=1}^{J_d} q_{gk d j}^{\beta_{gk d j}} \right\} \end{aligned}$$

From this distribution, we need to isolate the conditional distribution of each of the parameters. However, \mathbf{r}_m and \mathbf{z}_m are closely linked and separating them is somewhat difficult. Jointly, their conditional distribution is:

$$P(\mathbf{r}_m, \mathbf{z}_m | \mathbf{X}, \mathbf{R}, \mathbf{S}, \mathbf{Q}) \propto \prod_{g=1}^G R_g^{r_{mg}} \left(\prod_{k=1}^{C_g} (S_{gk} f(\mathbf{X}_m | \mathbf{Q}_{gk}))^{z_{mgk}} I \left(\sum_{k=1}^{C_g} z_{mgk} = r_{mg} \right) \right)$$

To find the conditional distribution for \mathbf{r}_m , we need to marginalize over \mathbf{z}_m by recognizing that:

$$\begin{aligned} P(r_{mg} = 1 | \mathbf{X}, \mathbf{R}, \mathbf{S}, \mathbf{Q}) &= \sum_{\mathbf{z}_{mg}} P(r_{mg} = 1, \mathbf{z}_{mg} | \mathbf{X}, \mathbf{R}, \mathbf{S}, \mathbf{Q}) \\ &\propto \sum_{\mathbf{z}_{mg}} R_g \prod_{k=1}^{C_g} (S_{gk} f(\mathbf{X}_m | \mathbf{Q}_{gk}))^{z_{mgk}} = R_g \sum_{k=1}^{C_g} S_{gk} f(\mathbf{X}_m | \mathbf{Q}_{gk}) \end{aligned}$$

Therefore, we can draw \mathbf{r}_m from:

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

Once we know which region the m th individual belongs to, we can draw \mathbf{z}_{mg} from^d:

$$(\mathbf{z}_{mg} | r_{mg} = 1, \mathbf{X}, \mathbf{S}, \mathbf{Q}) \sim \text{multinomial} \left(1, \left\{ \frac{S_{gk} f(\mathbf{X}_m | \mathbf{Q}_{gk})}{\sum_{k'=1}^{C_g} S_{gk'} f(\mathbf{X}_m | \mathbf{Q}_{gk'})} \right\}_{k=1,2,\dots,C_g} \right)$$

Next, given the regional identities, \mathbf{R} is drawn from:

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

Then, given the stock identities for each region, \mathbf{S}_g is drawn from:

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

Finally, for each stock within each region and for each locus, we generate \mathbf{Q}_{gkd} from:

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

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

^d This phrase is commented on in the section entitled "Technical Committee Review and Comments."

Simulations

Analyzing multiple simulated mixtures with Bayesian methods is somewhat challenging because no “canned” software is available to conduct automated analyses. For this reason, we were limited in the number of mixtures that could be analyzed. To simulate each fish, we randomly selected the stock of origin from the appropriate region, then, for each locus, we drew a genotype from the multinomial distribution using the observed baseline allele relative frequencies. We simulated 100 mixtures of 200 fish that were each composed of 100% Norton Sound chum, and analyzed them with a Western Alaska baseline. The baseline was composed of 53 SNPs and included 60 stocks representing 6 regions, including: Kotzebue Sound (5 stocks), Seward Peninsula (2 stocks), Norton Sound (12 stocks), Lower Yukon River (18 stocks), Kuskokwim River/Bay (17 stocks), and Bristol Bay (6 stocks). The mixtures were analyzed in 3 ways: 1) Pella-Masuda Model with the TFP, 2) Pella-Masuda Model with the RFP, and 3) RAM. The Pella-Masuda analyses were conducted in the R programming language utilizing the package BRUGS. The RAM analyses were also conducted within an R program, but the program called upon a C++ function that was developed at the Gene Conservation Laboratory to speed up analysis.^c For each mixture, one chain was run for 30,000 iterations, discarding the first 5,000 as burn-in. From the 25,000 iterations that were retained, posterior means of the stock proportions and the regional proportions were calculated. Also calculated were the means, central 90% quantiles, and root mean square errors of the 100 posterior means.

RESULTS

The mean and central 90% of the Norton Sound proportions were 0.831 (0.686–0.929) for the Pella Masuda Model TFP, 0.834 (0.696–0.932) for the Pella-Masuda Model RFP, and 0.880 (0.757–0.949) for the RAM (Table 1; Figures 3–5); and the root mean square errors were 0.091 for the TFP, 0.088 for the RFP, and 0.063 for the RAM (Table 1). For the Pella-Masuda Model, while both the TFP and the RFP showed very similar amounts of misallocation, the RFP tended to shift some of the misallocation away from the regions with the most stocks and into regions with fewer stocks (Figures 3–4). The RAM showed less misallocation than both prior structures of the Pella-Masuda Model in terms of point estimate and tightness of the central 90% quantiles, and tended to flatten out the amount of misallocation more evenly across the remaining regions (Figure 4).

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

Region	P-M TFP	P-M RFP	RAM
Kotzebue Sound	0.012	0.018	0.014
Seward Peninsula	0.004	0.011	0.01
Norton Sound	0.831	0.834	0.88
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

^c This sentence is commented on in the section entitled “Technical Committee Review and Comments.”

DISCUSSION

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

The rationale for why the RAM was expected to reduce the nonuniform bias can be seen by inspecting the regional identity probability:

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

This probability is a product of the regional contribution and a weighted average genotypic frequency, with the weights summing to one. Because the weights sum to one, the genetic component of this probability, i.e. the weighted average genotypic frequency, remains comparable regardless of the number of stocks within the region, which levels the playing field.^f The effect of this was seen in our simulation results. In our simulations, every mixture individual belonged to Norton Sound. Under the Pella-Masuda Model, when allocating the m th fish at each cycle, all 60 stocks competed for allocation of this fish. As can be seen in Figures 3 and 4, the larger regions were more successful at gaining this allocation simply because they have more stocks to compete with. However, under the RAM, when allocating the fish, only 6 regions were competing for allocation, each acting a single unit.

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

The RAM presented here is extended to only 2 levels of hierarchy of stocks within regions. However, it is conceivable to expand this model to further levels of hierarchy, such as sub-stocks within stocks, and stocks within regions. Such a model may be useful in situations where multiple levels of structure exist.

^f This sentence is commented on in the section entitled "Technical Committee Review and Comments."

ACKNOWLEDGEMENTS

The Technical Document series served as a record of communication between the Alaska Department of Fish and Game Commercial Fisheries Division and the Western Alaska Salmon Stock Identification Program (WASSIP) Technical Committee during the implementation of the program. The authors would like to thank the WASSIP Technical Committee and Advisory Panel for their constructive input on each of the documents throughout the project. The authors would also like to thank Erica Chenoweth who coordinated and prepared the Technical Document series for publication and Publication Specialists Amy Carroll and Joanne MacClellan for implementing the series into Regional Information Reports.

REFERENCES CITED

- Dann, T. H., C. Habicht, J. R. Jasper, H. A. Hoyt, A. W. Barclay, W. D. Templin, T. T. Baker, F. W. West, and L. F. Fair. 2009. Genetic stock composition of the commercial harvest of sockeye salmon in Bristol Bay, Alaska, 2006-2008. Fishery Manuscript Series No. 09-06.
- Fournier, D. A., T. D. Beacham, B. E. Ridell, and C. A. Busack. 1984. Estimating stock composition in mixed stock fisheries using morphometric, meristic, and electrophoretic characteristics. *Canadian Journal of Fisheries and Aquatic Sciences*. 52:1688–1702.
- Gelman, A., and D. Rubin. 1992. Inferences from iterative simulation using multiple sequences. *Statistical Science*. 7:457–511.
- Okuyama, T. and B. M. Bolker. 2005. Combining genetic and ecological data to estimate sea turtle origins. *Ecological Applications*. 15(1):315–325.
- Pella, J. and M. Masuda. 2001. Bayesian methods for analysis of stock mixtures from genetic characters. *Fishery Bulletin*. 99:151–167.
- Pella, J. J., and G. B. Milner. 1987. Use of genetic marks in stock composition analysis. Pages 247-276 [In] N. Ryman and F. Utter, editors. *Population Genetics and Fishery Management*. Washington Sea Grant Program WASHU-B-87-001.
- Rannala, B. and J. L. Mountain. 1997. Detecting immigration by using multilocus genotypes. *Proceedings of the National Academy of Sciences*. 94: 9197–9201.
- Wood, C. C., S. McKinnell, T. J. Mulligan, and D. A. Fournier. 1987. Stock identification with the maximum-likelihood mixture model: sensitivity analysis and application to complex problems. *Canadian Journal of Fisheries and Aquatic Sciences*. 44:866–881.

TECHNICAL COMMITTEE REVIEW AND COMMENTS

WASSIP Technical Document 7 Regional Allocation Model (RAM)

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

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

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

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

- Stocks are organized hierarchically into 2 or more regions or Reporting Groups (RGs).
- The RGs have the same number of actual populations but different numbers of populations that have been sampled for the baseline.
- A flat prior of stock contribution is computed as $1/n$, where n is the total number of populations in the baseline.
- In this scenario, the RGs that have the most populations in the baseline will tend to attract the most spurious contribution assigned to low- or noncontributing stocks.

The solution to this problem proposed by Technical Document 7 is two-fold:

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

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

Consider the following scenario:

- Stocks are organized hierarchically into 2 or more regions or RGs.
- The RGs have different numbers of actual populations, and each actual population has been sampled for the baseline.
- Each population has the same size and productivity.

Under this scenario, the appropriate priors for each RG are proportional to the number of stocks in the baseline, and enforcing equal RG priors as in item 1) above could be expected to reduce accuracy of the estimates.

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

For the particular case of separating stocks in mixtures taken from the WASSIP study area, the authors might think about the potential for using semi-informative priors, and investigate whether the priors have an appreciable effect on the results. For example, abundance varies greatly among the stocks/regions investigated; proximity of these stocks to the WASSIP area varies as well, and there is some rudimentary oceanic distribution information from tagging studies. Hopefully, the results aren't too sensitive to the priors on stock composition, but if they are, these priors should receive careful attention. In case of sensitivity, priors should be chosen based on the best biological information and possibly partially on management priorities. The effects of priors on estimates for small stocks should get particularly careful consideration. If the priors weight each region equally, and some of these small stocks get treated like a region, the priors could potentially dominate the results and strongly overweight their contributions.

Comments keyed to specific lines:

Page 2, 1st ¶, third sentence (note a): this is true only if some method has been used to account for unsampled alleles

Page 2, 4th ¶, first sentence (note b): isn't this a null hypothesis rather than an assertion?

Page 5, antepenultimate ¶, last sentence (note c): is ragged matrix a real term?

Page 7, middle (note d): "once we know ..." ... do you mean, "once we have estimated"?

Page 8, 1st ¶, eighth sentence (note e): what exactly did the C++ routine do?

Page 9, near middle (note f): we agree that in the example chosen, the new method helps to "level the playing field." However, as discussed above, forcing equal RG priors is not a sound general strategy for leveling the playing field.

Figure 1 (note g): how was the individual stock of origin for each Norton Sound fish in the simulated mixtures chosen?

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

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

ADDITIONAL COMMENTARY

These comments are excerpted from Jerry Pella's email message sent on April 14, 2010 in reply to an email from James R. Jasper on the Regional Allocation Model (RAM).

I read your draft and compared it with some earlier work by Toshinori Okuyama and Benjamin Bolker 2005. Combining genetic and ecological data to estimate sea turtle origins. *Ecological Applications* 15(1): 315-325. Their appendices are also applicable and available at the Ecological Archives (A015-009-A1 and A2) site. Although these authors only considered mitochondrial DNA haplotypes, the remainder of their model appears identical to yours with their "gyres" equivalent to your "regions" and their "rookeries" equivalent to your "stocks". They had information about the magnitudes of the rookeries and I suspect their motivation in developing a gyres model was to allow differing covariate relationships for magnitudes and contributions among gyres.

Your note begins with the statement that an inherent problem of nonuniform bias occurs due to differing numbers of stocks in regions. I have not been aware of this bias, not having heard the complaint previously. I can understand that by beginning the mcmc computations from a point of equal stock contributions, a region with many stocks could receive a too-high allocation of mixture individuals before convergence of the mcmc chain, but as convergence is approached, that kind of "bias" should decrease and eventually disappear. On the other hand, if convergence of the chain does not occur, regions with many stocks could well be allocated too many mixture individuals. Your solution to overcoming this unequal regional allocation involves first setting a prior on the regional contributions and second, on the stock contributions within the regions. If, for example, you set a low information prior on the regions (equal regional proportions), the effect as compared to the simple Pella & Masuda model should be to reduce the numbers of mixture individuals allocated to regions with many stocks and to increase the numbers of mixture individuals allocated to regions with few stocks at the beginning of the chain of computations, but eventually this difference from the Pella & Masuda computations should disappear as convergence of the chain occurs. If the difference does not disappear, then the regional prior must be informative and cannot be viewed as low information. If so, the regional prior would have to be chosen with some substantive basis. However, until shown otherwise, I would expect both mcmc chains to converge to essentially the same regional and individual stock posterior distribution for mixture proportions if the regional and stocks with region contribution priors were thought to be of low information.

The regional allocation of mixture individuals is in proportion to regional weighted averages for relative frequencies of mixture genotypes with weights equal to stock proportions within regions. If stocks vary genetically a lot within regions, a regional average may not represent the genotype relative frequencies of the individual stocks within the region very well and the allocation may perform poorly. On the other hand, if the variation among regions is large relative to within regions, it could work well and speed convergence. You suggest in the draft that a lot of overlap among stocks between regions occurs, which sounds like regional averages might not work well to me.

FIGURES

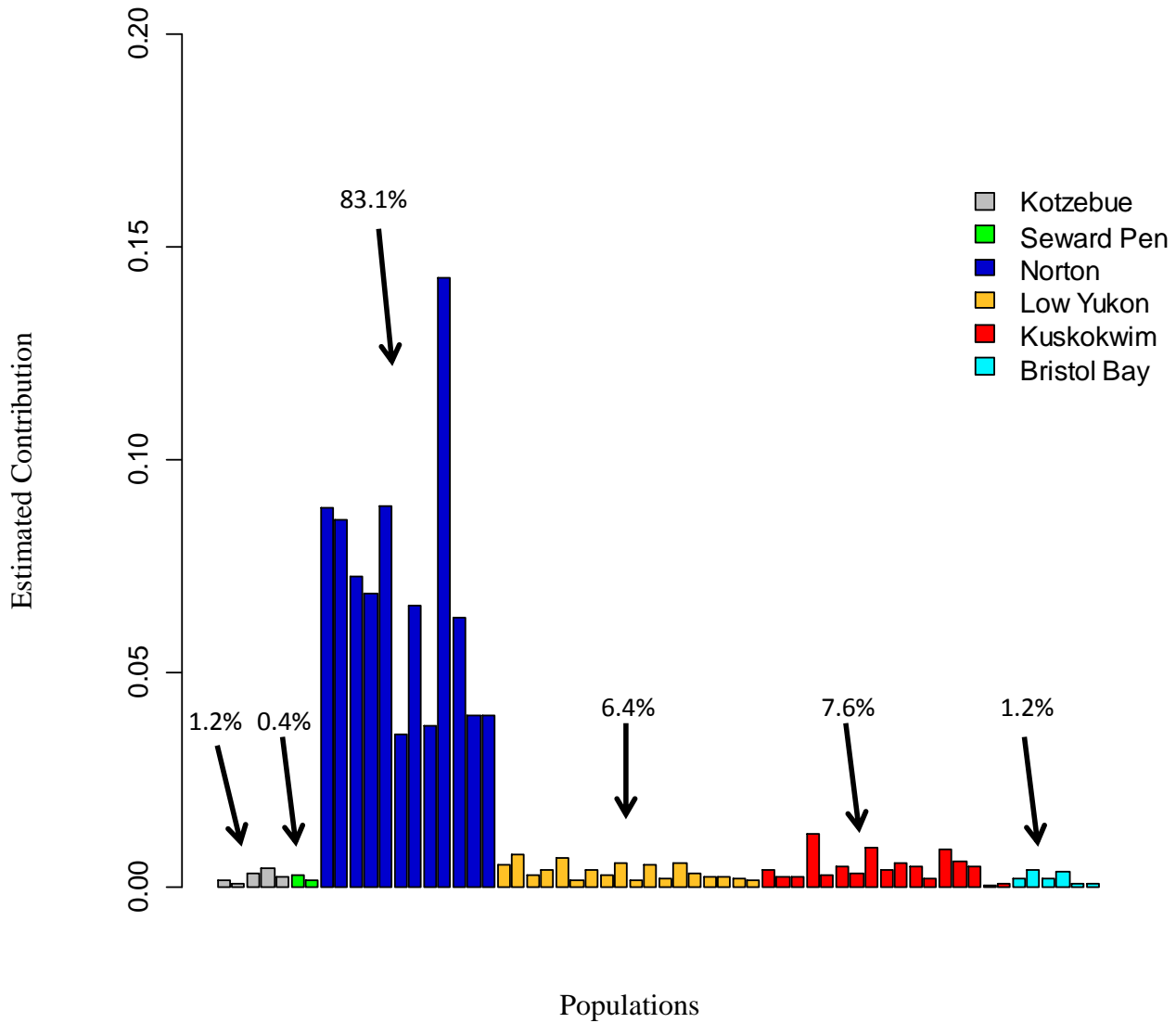


Figure 1.—Simulation results for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model shown at the individual stock level. The height of the bars represents the mean of 100 repetitions. An equal prior “count” of one divided by the number of stocks was given to each stock.^g Regional means are indicated.

^g This figure is commented on in the section entitled “Technical Committee Review and Comments.”

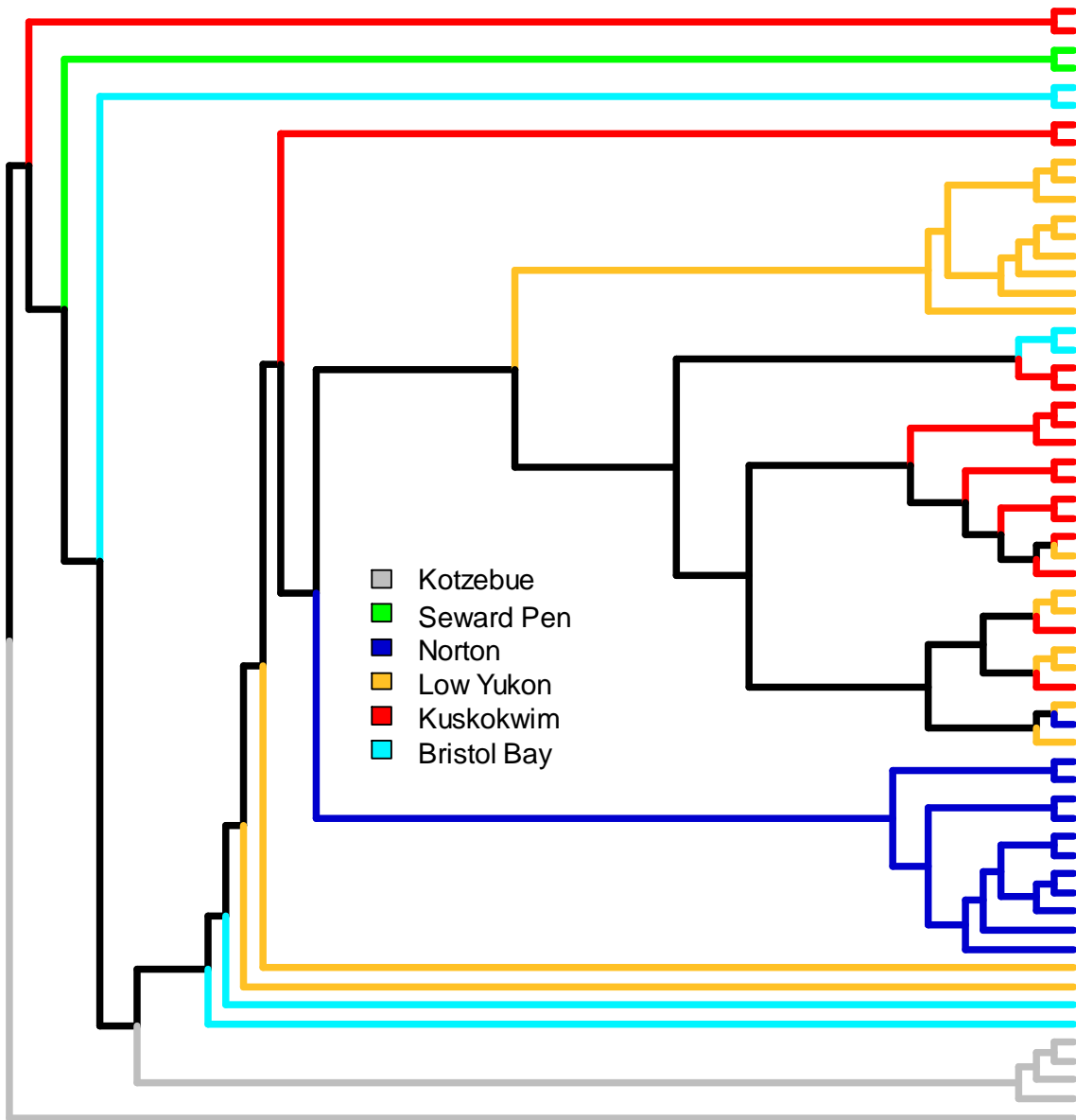


Figure 2.—Unweighted pair-group method (UPGMA) tree of pair-wise F_{ST} for 60 stocks of Western Alaska chum demonstrating that Norton Sound chum are more genetically similar to Lower Yukon and Kuskokwim than the other regions.

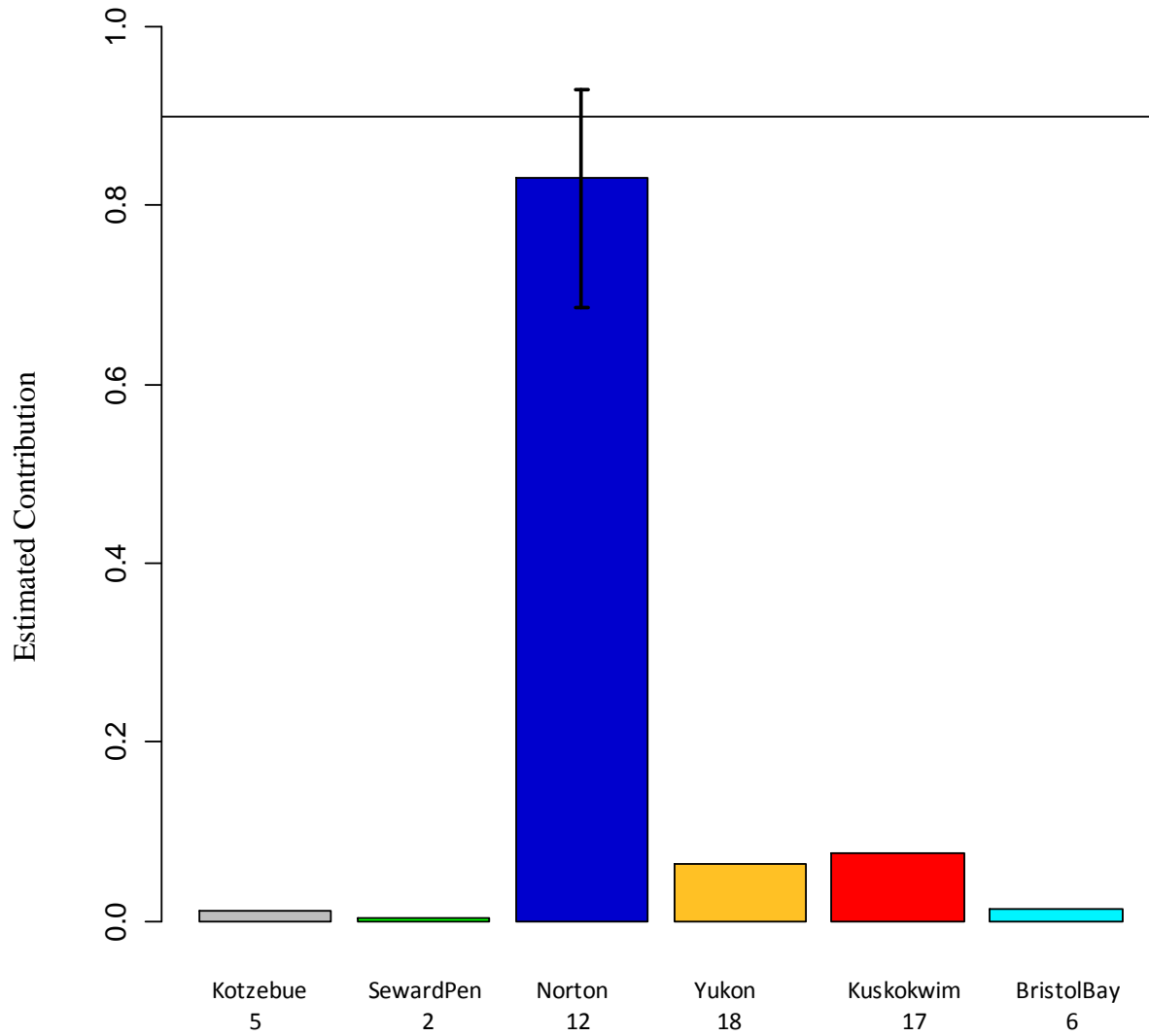


Figure 3.—Simulation results for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model using the True Flat Prior. The height of the bars represents the mean of 100 repetitions with the 90% credibility interval indicated. The horizontal rule is 90% correct allocation. Numbers under labels are the number of stocks within the region. These results are the same as shown in Figure 1 with the stock proportions summed into regions.

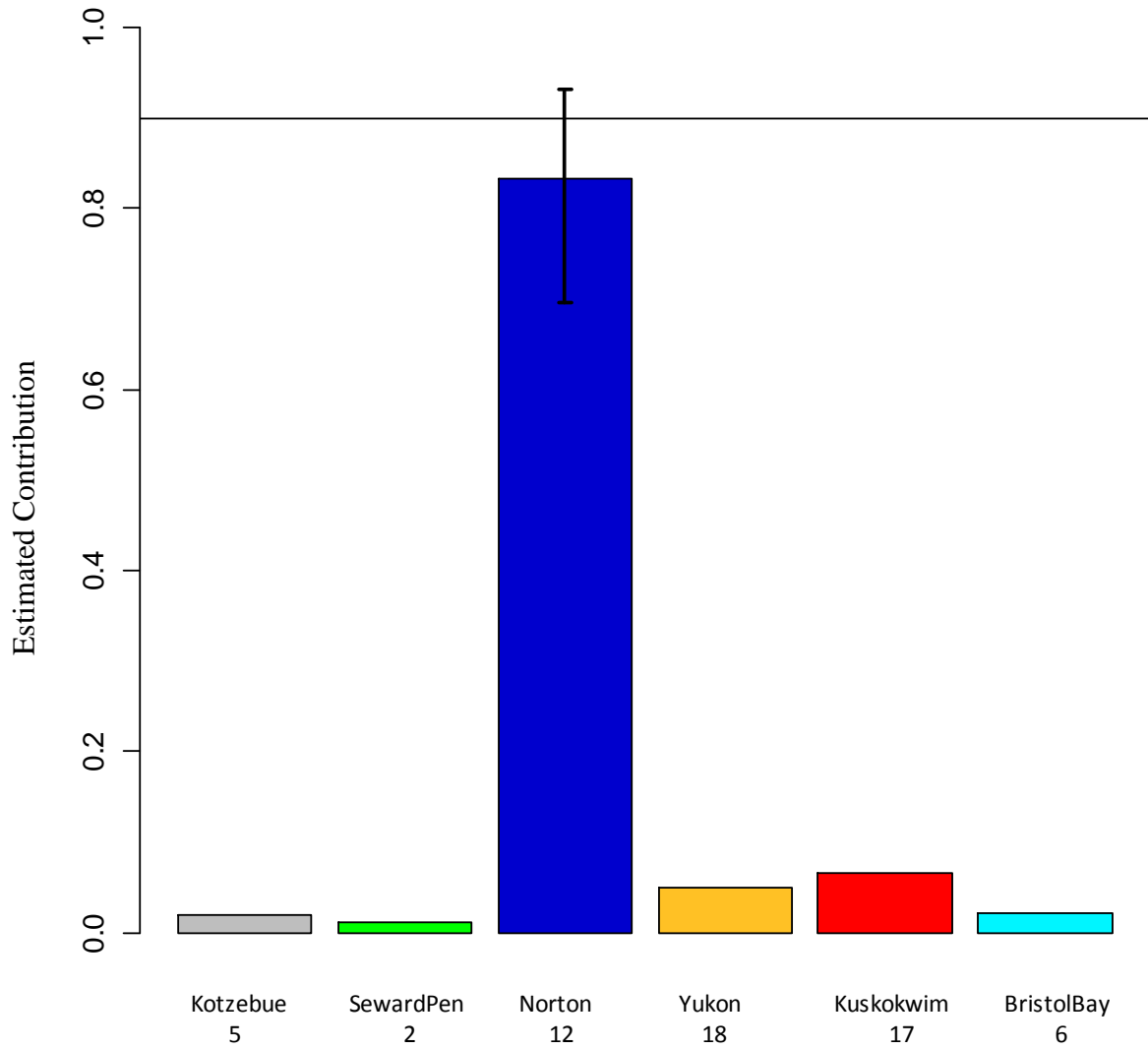


Figure 4.— Simulation results for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model using the True Flat Prior. The height of the bars represents the mean of 100 repetitions with the 90% credibility interval indicated. The horizontal rule is 90% correct allocation. Numbers under labels are the number of stocks within the region. Numbers under labels are the number of stocks within the region.

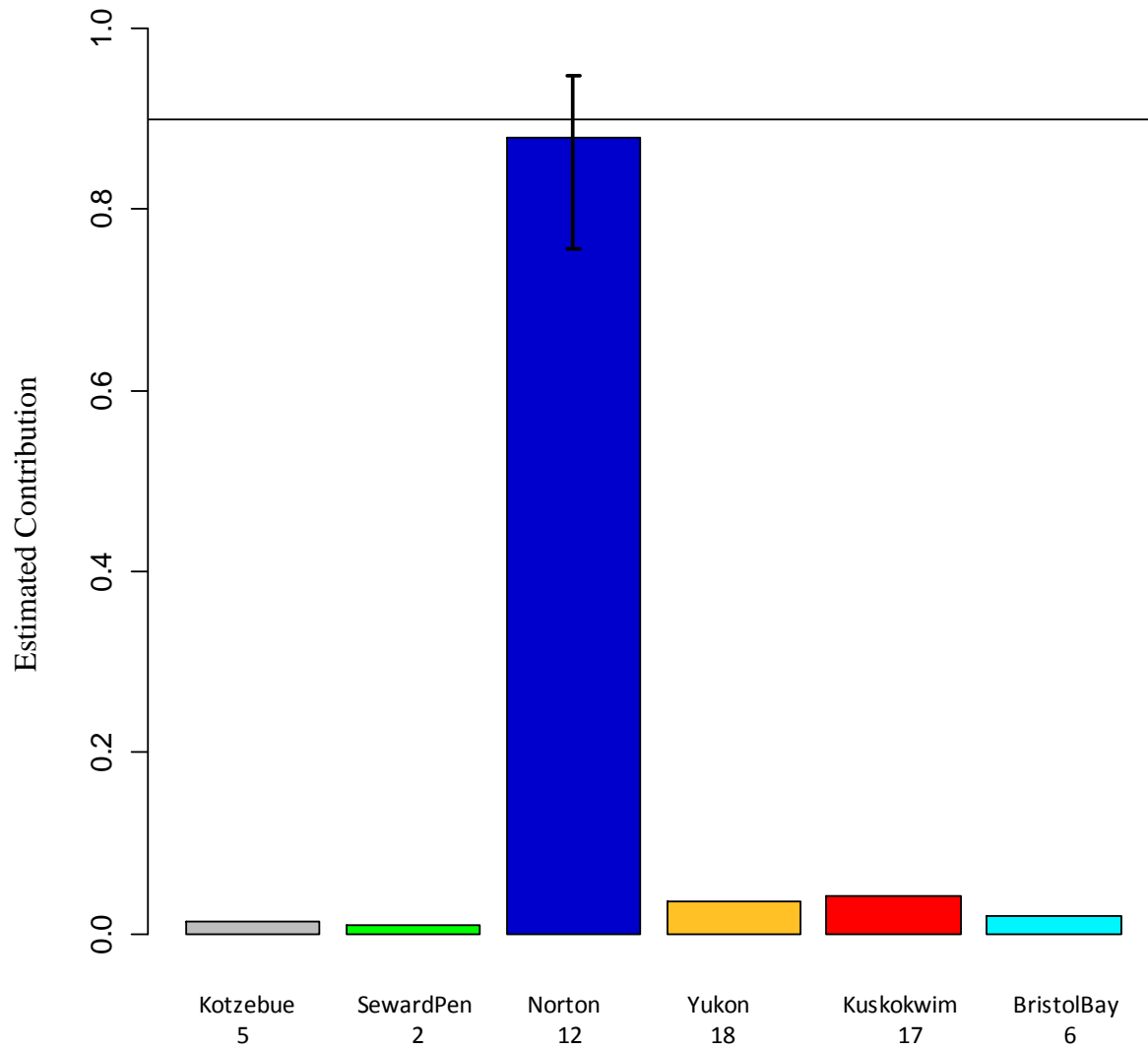


Figure 5.— Simulation results for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model using the True Flat Prior. The height of the bars represents the mean of 100 repetitions with the 90% credibility interval indicated. The horizontal rule is 90% correct allocation. Numbers under labels are the number of stocks within the region.