

1
2
3
4
5

Version: 1.0

Title: Chum salmon SNP selection process outline

Version:
1.0

Authors: N. DeCovich, J. Jasper, C. Habicht, W. Templin

Date: December 14, 2010

6 **Introduction:**

7 Early in the development process for the Western Alaska Salmon Stock Identification Project
8 (WASSIP) it was clear that the resolution possible for chum salmon spawning in western Alaska
9 regional areas (Norton Sound, lower Yukon and Kuskokwim rivers, and Bristol Bay) was not
10 going to be sufficient to meet the standards set by the Advisory Panel (AP) with available genetic
11 markers, including the recently developed SNP markers (see Technical Document 4 for the
12 current panel of 53 SNPs). These four regional areas define important units for management, yet
13 when treated as separate reporting groups each performed below the 90%-correct-allocation level
14 using the 53-marker set. The Department began the process of discovering additional SNP
15 markers for chum salmon through a contract with the International Program for Salmon
16 Ecological Genetics (IPSEG; <http://www.fish.washington.edu/research/ipseg/research.html>) at
17 the University of Washington. These efforts were based on cDNA sequences from two chum
18 salmon sampled from the Susitna and Delta rivers. This process has been described in a
19 manuscript that has been published in *Molecular Ecology Resources* (Seeb et al. 2011) which is
20 provided as Technical Document 9. This process added 37 validated SNPs to those already
21 available for chum salmon for use in WASSIP. Subsequent rounds of SNP development at the
22 University of Washington were based on 16 fish from four populations from Western Alaska and
23 increased the total number of described SNPs to 228 (Grau et al. in prep).

¹ 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.

24 Here we describe the process that we intend to use to select the set of 96 SNPs that maximizes
25 the likelihood of providing the resolution necessary to meet the objectives of WASSIP. A
26 similar process was recently completed with the selection of 96 SNP markers for use with
27 sockeye salmon and is described in Technical Document 6, “Selection of the 96-SNP marker set
28 for sockeye salmon.” However, the selection of chum salmon SNPs will be significantly
29 different from that used for sockeye salmon. There are many more SNPs available for chum
30 salmon than were available for sockeye salmon (124 SNPs), and more emphasis is placed on
31 selecting markers to distinguish among regional areas (Norton Sound, Yukon summer,
32 Kuskokwim summer, Western Bristol Bay, and Eastern Bristol Bay) within coastal western
33 Alaska (CWAK).

34 **Method:**

35 I. *Pre-ADF&G selection:* Markers were developed under contract at the IPSEG laboratory.

- 36 a. TaqMan assays were developed or available for a total of 228 SNPs including the
37 original 53 SNPs.
- 38 b. Markers were assayed in 80 - 96 individuals from each of 30 populations (Table
39 1; Figure 1) chosen from across the species range. Ten of these populations were
40 from CWAK (Figure 2).
- 41 c. Of the 228 markers surveyed, 188 markers have been determined to perform
42 adequately in the laboratory and have a reasonable level of variation. Only these
43 markers will be passed on from IPSEG to ADF&G for further analysis.

44 II. *Unranked measures:* The measures in this section will be given veto power. Markers will be
45 discarded if they do not pass the following tests.

- 46 1. Hardy-Weinberg Equilibrium. Conformance to HWE will be measured using the
47 program Genetic Data Analysis (GDA; Lewis & Zaykin 2001). GDA uses the
48 methods described in Genetic Data Analysis II (Weir 1996). Markers out of HWE at
49 $\alpha = 0.05$ in more than 5 populations or out of HWE at $\alpha = 0.001$ in more than one
50 population will be dropped.
- 51 2. Linkage Disequilibrium/Phase. Linkage Disequilibrium will be measured with the
52 program GDA.

- 53 a. Significant disequilibrium between markers will be determined using the
54 sequential Bonferroni with an overall level $\alpha=0.05$ for each marker set adjusted by
55 the number of populations.
- 56 b. For marker sets that exhibit disequilibrium, we will next determine whether
57 combining linked markers or discarding a marker is most useful for MSA. To do
58 this with a pair of linked markers we will set up three treatment files:
- 59 i. Marker A combined with marker B (“composite phenotype”; Habicht et al.
60 2010);
- 61 ii. Marker A retained and marker B excluded; and
- 62 iii. Marker B retained and marker A excluded.

63 This can be extended to larger linked groups if necessary. We will use f_{orca}
64 (Rosenberg 2005) and measure correct individual assignment to population for the
65 three treatments. The treatment with the best average correct assignment will be
66 selected for further analyses. This method is similar to the methods outlined in
67 Ackerman et al. (In press) where GENECLASS (Piry et al. 2004) was used for the
68 assignment software.

69 III. *Ranked or scored measures of population structure and MSA performance*: The measures in
70 this phase of the selection process are either ranked or scored and then weighted. Highest
71 weighting is given to measures associated with variation among CWAK populations.
72 Weights are given as percentages and sum to 100%.

- 73 1. CWAK –specific measures [84% of total].

74 Question addressed: What are the best markers for distinguishing among populations
75 or regions within CWAK? This is the most difficult portion of the range to
76 distinguish population structure, yet resolution within this area is central to the
77 objectives of WASSIP.

- 78 a. Among populations (24%)

- 79 i. Overall F_{ST} among the 10 CWAK populations. The F_{ST} values calculated
80 from individual markers will be linearly scaled between 0.0 (lowest) and 1.0
81 (highest) and used as scores.

82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110

- b. Among regions (60%)
- i. Overall θ_P among the 5 CWAK regions. θ_P for each marker will be calculated via a three-level hierarchical ANOVA (Weir, 1995), in which populations from CWAK are organized into five regions (Table 1; Figure 2). The θ_P values calculated from individual markers will be linearly scaled between 0.0 and 1.0 and used as scores. (See Figure 2; 15%)
 - ii. f_{orca} (Rosenberg 2005) with backward elimination marker selection algorithm method using the five CWAK regions as reporting groups. This method is similar to BELS (Bromaghin 2008) in that it starts with all markers and then sequentially eliminates the marker that provide the least amount of regional discrimination (Technical Document 10). Each marker is then ranked according to the order in which they were eliminated. To then score each marker, we sequentially add markers according to their rank, starting with the most informative marker, and calculate f_{orca} at each step. The resulting f_{orca} values can then be linearly scaled between 0.0 and 1.0, with one corresponding to the most informative marker. BELS is too time-consuming to be used and relies on a simulation method that may introduce bias. (30%)
 - iii. $\theta_{S(P)} = \theta_S - \theta_P$ for population pairs from adjacent CWAK regions. $\theta_{S(P)}$ for each marker will be calculated via a three-level hierarchical ANOVA, in which populations from adjacent regions are paired. This quantity is a measure of the differentiations among populations within pairs. The four population pairs from adjacent regions with smallest pairwise F_{ST} will be chosen for these tests. The $\theta_{S(P)}$ values calculated from individual markers will be linearly scaled between 0.0 and 1.0 and used as scores. (15%)

2. Pacific-wide measures [10% of total].

Question addressed: What are the best markers for distinguishing among large-scale regions across the species range? Some of the WASSIP fisheries are known to intercept chum salmon from both the western and southeastern extent of the range.

111 These measures will ensure that broad-scale regions will be identifiable in WASSIP
112 fishery samples.

- 113 a. Principle Component Analysis. The amount of variation explained by each
114 marker will be linearly scaled between 0.0 and 1.0 and used as scores.
- 115 i. The amount of variation associated with each marker in the first principle
116 component (3%)
 - 117 ii. The amount of variation associated with each marker in the second
118 principle component (3%)
 - 119 iii. The amount of variation associated with each marker over the set of
120 principal components that explain 80% of the variation (4%)
- 121 3. Outside Alaska, regional measures [6% of total].

122 Question addressed: What are the best markers for distinguishing between population
123 pairs within or between certain regions outside of Alaska? This is expected to
124 provide insight into markers important for distinguishing broad-scale population
125 structure and is considered to insure a useable panel of SNPs for research groups
126 outside of Alaska. (See Figure 3)

- 127 a. Within Japan. Calculate the F_{ST} between populations selected from Honshu and
128 Hokkaido islands (2%). The F_{ST} values calculated from individual markers will
129 be linearly scaled between 0.0 and 1.0 and used as scores.
- 130 b. Between Southeast Alaska and Northern British Columbia. Calculate the F_{ST}
131 between population pairs selected from Southeast Alaska and Northern British
132 Columbia (2%). The F_{ST} values calculated from individual markers will be
133 linearly scaled between 0.0 and 1.0 and used as scores.
- 134 c. Between Southern British Columbia and Washington. Calculate the F_{ST} between
135 population pairs selected from Southern British Columbia and Washington (2%).
136 The F_{ST} values calculated from individual markers will be linearly scaled between
137 0.0 and 1.0 and used as scores.

138 IV. *Final considerations:* The candidate SNPs will be ordered from best to worst with respect to
139 the measures in Section III above. The measures in this section (IV) will be performed on
140 the top 96 candidates based on the measures in Section III (above). If a marker is discarded

141 due to laboratory performance, the next highest-rated marker from Section III will be
142 evaluated.

- 143 1. Performance at the IPSEG Laboratory. Assay performance will be evaluated on three
144 criteria. High-ranking markers that have poor laboratory performance and lower-
145 ranked markers that are difficult to score will be dropped and replaced with the next
146 highest-ranking marker. The process will continue until 96 markers are selected. We
147 incorporate laboratory performance here to avoid the need to examine assay
148 performance of markers that provide little useful MSA performance. Laboratory
149 performance will be evaluated with the following measures as used in the sockeye
150 selection process (Technical Document 6):
 - 151 a. Cluster tightness (See Figure 4)
 - 152 b. Cluster alignment (See Figure 5)
 - 153 c. Drop-out rates (See Figure 6)
- 154 2. Final evaluation using simulations to test for loss of MSA resolution for
155 distinguishable regions generally outlined in Seeb et al. (2011). Simulations will be
156 conducted using the selected markers to ensure that the reporting groups represented
157 in this data set that were distinguishable in Seeb et al. (2011) continue to be
158 distinguishable (> 90% correct allocation). Matching exact reporting groups will not
159 be possible, but reasonable approximations will be tested. These reporting groups will
160 include (corresponding population numbers from Table 1 in parentheses): Japan (1,2),
161 Russia (3,4), Kotzebue Sound (5,6), CWAK (7,8,9,10,13,14,15,16), Yukon Fall
162 (11,12), Eastern Bristol Bay (17,18), North Alaska Peninsula (19,20), South Alaska
163 Peninsula (21,22), Southcentral Alaska (23,24), Southeast Alaska/BC (25,26,27,28),
164 and Washington (29,30). Mean correct allocations in the Seeb et al. (2011) study
165 ranged from 85% to 99%, with the majority of reporting regions allocating above
166 90%. The results from our analysis are expected to be optimistic given that regions
167 are represented by only a few populations. Therefore, mean correct allocations to
168 reporting groups below 90% will trigger addition of markers that were highly ranked
169 from sections III.2 and III.3. As markers are added, the lowest-ranked markers from
170 the III.1 process will be dropped. Markers will be added and dropped following these
171 rules until the resolution to these broader reporting groups exceeds 90%.

172 3. Laboratory performance in ADF&G. All 188 SNPs will be assayed in the Gene
173 Conservation Laboratory on 3,032 chum salmon originating from Prince William
174 Sound as part of a Pacific Coast Salmon Recovery Fund project. This will allow us to
175 confirm assay performance in our lab.

176 **Literature Cited:**

177 Ackerman, M. W., C. Habicht, and L. W. Seeb. In Press. SNPs under diversifying selection
178 provide increased accuracy and precision in mixed stock analyses of sockeye salmon
179 from Copper River, Alaska and nearby coastal areas. Transactions of the American
180 Fisheries Society. XX:XXX-XXX

181 Bromaghin, J. F. 2008. BELS: backward elimination locus selection for studies of mixture
182 composition or individual assignment. Molecular Ecology Resources 8: 568-571.

183 Habicht, C., L. W. Seeb, K. W. Myers, E. V. Farley, and J. E. Seeb. 2010. Summer-fall
184 distribution of stocks of immature sockeye salmon in the Bering Sea as revealed by
185 single-nucleotide polymorphisms. Transactions of the American Fisheries Society
186 139(4):1171-1191.

187 Lewis P.O., and D. Zaykin. 2001. GENETIC DATA ANALYSIS: computer program for the
188 analysis of allelic data, version 1.0 (d16c) Free program distributed by the authors from
189 <http://lewis.eeb.uconn.edu/lewishome/software.html>.

190 Piry, S., A. Alapetite, J. M. Cornuet, D. Paetkau, L. Baudouin, and A. Estoup. 2004.
191 GENECLASS2: A software for genetic assignment and first-generation migrant
192 detection. Journal of Heredity 95(6):536-539.

193 Rosenberg, N.A. 2005. Algorithms for selecting informative marker panels for population
194 assignment. Journal of Computational Biology 12(9):1183-1201.

195 Seeb, J.E., C.E. Pascal, E.D. Grau, L.W. Seeb, W.D. Templin, S.B. Roberts, and T. Harkins.
196 2011. Transcriptome sequencing and high-resolution melt analysis advance SNP
197 discovery in duplicated salmonids. Molecular Ecology Resources doi: 10.1111/j.1755-
198 0998.2010.02936.x.

199 Seeb, L.W., W.D. Templin, S. Sato, S. Abe, K.I. Warheit, and J.E. Seeb. 2011. Single nucleotide
200 polymorphisms across a species' range: implications for conservation studies of Pacific
201 salmon. Molecular Ecology Resources xxx: xx-xx

202 Weir, B. S. 1996. Genetic Data Analysis II. Sinauer Associates, Inc., Sunderland, Mass.

203

204 **Specific questions for the Technical Committee:**

- 205 1. Is our approach to linkage disequilibrium and HWE reasonable?
206 2. Is our method to determine the relative value of different treatments of linked markers
207 advisable? Is the use of f_{orca} as a measure appropriate?
208 3. Are the tests appropriately structured to provide a set of SNPs that will perform well for
209 WASSIP?
210 4. Does the weighting applied to each set of tests seem reasonable?
211 5. Are there other measures that would be more appropriate?

212

213 General comments: The approach proposed here borrows useful ideas from the approach used
214 for sockeye salmon (described in Technical Document 6) but appears to be more streamlined and
215 efficient. The text is a bit confusing about how the laboratory screening will occur. At line 41,
216 the report states, "Of the 228 markers surveyed, 188 markers have been determined to perform
217 adequately in the laboratory and have a reasonable level of variation. Only these markers will be
218 passed on from IPSEG to ADF&G for further analysis." This implies that data quality issues in
219 the laboratory have already been evaluated prior to screening loci for power to discriminate
220 populations. However, at line 140 another process is described that seems to involve iterative
221 consideration of discriminatory power and laboratory performance.

222

223 Responses to specific questions:

224 1. *Is our approach to linkage disequilibrium and HWE reasonable?*

225 For the most part, but we have several comments to consider.

- 226 1) For both types of analyses, it is important to ensure that the baseline populations represent
227 single panmictic populations. If not, a Wahlund effect could cause both HW and LD
228 departures that appear to be data quality issues but actually reflect population mixture.
229 2) For both types of analyses, be careful about only using results of tests of statistical
230 significance. You are really interested in the magnitude of the effect size here, but P values
231 also depend heavily on sample sizes. Also, the direction of departure (e.g., heterozygotes
232 excess or deficiency) can be informative about potential causes.
233 3) The LD analyses will consider pairs of loci, of which there are $n(n-1)/2$ possible comparisons
234 for n loci. Since n could be 200 or more, this represents a huge number of pairwise
235 comparisons, each of which could be conducted for many different populations. Using the
236 Bonferroni correction here would require consideration of tiny P values, which could lead to
237 unpredictable results. It is probably more useful to screen for pairs of loci that are
238 consistently out of equilibrium (using the nominal alpha level) in multiple populations. Some
239 consideration of effect size (the magnitude of LD) would also be useful in evaluating how
240 serious a problem any deviations are likely to cause.

241

242 2. *Is our method to determine the relative value of different treatments of linked markers*
243 *advisable? Is the use of FORCA as a measure appropriate?*

244 The general procedure described at lines 56-68 of Document 8 seems reasonable, as does
245 the logic for using a procedure that assigns entire individuals rather than making fractional

246 assignments. With the caveats noted below, *fORCA* should be ok as a means to assess *relative*
247 power for correct assignment.

248

249 3. *Are the tests appropriately structured to provide a set of SNPs that will perform well for*
250 *WASSIP?*

251 The proposed methods should produce a set of SNPs with high power to resolve stock
252 identification problems in Western Alaska.

253

254 4. *Does the weighting applied to each set of tests seem reasonable?*

255 The weights chosen are obviously somewhat arbitrary but do not appear to be unreasonable.
256 Because of the applied focus of this project, it is appropriate to assign greater weight to markers that have
257 high power for the local areas of interest. However, we were pleased to see that the criteria include non-
258 trivial weight to markers with wider geographic relevance (10% weight for Pacific Rim individual
259 populations, plus 6% for major non-Alaska groups). This will help ensure that the considerable efforts
260 here to develop markers will have much broader application to the scientific and fishery management
261 communities.

262

263 Minor comments:

264 In the proposed PCA analysis for Pacific-wide assessments, part (iii) is partially redundant as it
265 will include information already used for (i) and (ii)

266 Outside Alaska: we don't necessarily disagree with the particular comparisons proposed, but the
267 rationale for choosing them is not given.

268

269 5. *Are there other measures that would be more appropriate?*

270 Can't think of any offhand.

271

272 General comments about bias and *fORCA*

273 It is important to distinguish between two different types of biases that can potentially arise in
274 evaluations such as those proposed here.

275 The first type of bias, described by Anderson et al. (2008), occurs when one is interested in
276 assessing the power of a particular set of markers to resolve the composition of a mixture comprised of
277 individuals from a specified group of source populations. The ideal way to do this is to create simulated
278 mixtures of individuals, with the genotype of each individual being chosen based on actual allele
279 frequencies in one of the (randomly chosen) source populations. The bias arises because we never know
280 the actual allele frequencies—we only have samples. Because of random sampling error, allele
281 frequencies in samples from the baseline populations will on average be more divergent than are the true
282 population allele frequencies. On average, this factor inflates F_{st} among baseline samples by the
283 magnitude $1/(2S)$, where S is the baseline sample size. When simulated mixtures are constructed using

284 these baseline allele frequencies (which appear more different than the populations actually are), the
285 population assignments will tend to be overly optimistic. Furthermore, the relative importance of
286 sampling error (and hence the bias) will be larger when true genetic differences among populations are
287 very small—as occurs with Western Alaska chum salmon. Anderson et al. (2008) described a simple
288 leave-one-out procedure that eliminates the bias, but the routine described at lines 41-50 of Document 10
289 would be subject to this type of bias.

290 The second type of bias, described by Anderson (2010), applies to locus-selection programs. The
291 bias is not in the locus selection *per se*, but rather in the evaluation of power of the resulting set of loci for
292 population assignment. Anderson (2010) showed that the bias arises because none of the commonly-used
293 software programs for locus selection (including BELS) use proper cross validation. Instead, some of the
294 information used to select the panel of loci is also used to evaluate its performance, and this leads to an
295 overly optimistic assessment of assignment power. We did not see any indication that the combined
296 *fORCA*-BELS approach proposed in Document 10 would *not* be subject to this type of bias. Also,
297 although the authors list 4 methods Rosenberg (2005) evaluated for selecting subsets of loci, they don't
298 explain why they did not consider any of them for the current project.

299 One reason that proper cross-validation is often not done is that it is costly in terms of
300 information content. The “gold standard” of cross validation is to split the data in half: the first half is
301 used to develop the algorithm, the second half to evaluate its performance. However, doing this means
302 that the algorithm is likely to be less precise because it is based on less data. Researchers are thus
303 typically faced with a trade-off between precision in developing the best algorithm (use all the data in the
304 first step) and the downstream consequences (subsequent assessments of performance using the same data
305 will tend to be overly optimistic). Anderson (2010) suggested a simple modification to the cross-
306 validation procedure that retains most of the information without leading to appreciable bias in assessing
307 performance.

308 In summary, both types of biases can lead to overly optimistic assessments of power, which
309 should be a concern given the stated goals of the project. For applications that only consider relative
310 power, these biases might not be important. Also, it might be the case that the proposed locus-selection
311 approach is perfectly fine for selecting an optimal panel of loci, but that the estimates of power to be
312 expected when that panel is applied to real data are biased upwards.

313 Text at lines 84-91 of Document 10 seems to acknowledge at least the bias problem identified by
314 Anderson et al. (2008), but it is not clear that both of the potential sources of bias described above have
315 been fully considered in the documents we reviewed. This topic merits closer scrutiny to determine the
316 optimal way to proceed given project goals.

317

318 Anderson, E.C., R.S. Waples, S.T. Kalinowski. 2008. An improved method for estimating the accuracy of
319 genetic stock identification. *Canadian Journal of Fisheries and Aquatic Sciences* 65:1475-1486.

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

322

323 Table 1. Population set used in this analysis. Map numbers correspond to numbers in Figure 1.

ADF&G region	Population	Sample size	Map Number
Japan	Tokachi River	80	1
	Gakko River late	80	2
Russia	Amur River summer	95	3
	Palana River	95	4
Kotzebue Sound	Kiana River	95	5
	Inmachuk River	95	6
¹ Norton Sound	Kwiniuk River	95	7
	Unalakleet River	95	8
¹ Yukon summer	Andreafsky River - East Fork weir	95	9
	Nulato River	95	10
Yukon fall	Fishing Branch	95	11
	Kluane River	95	12
¹ Kuskokwim summer	Salmon River	95	13
	Kanektok River weir	95	14
¹ Western Bristol Bay	Osviak River	95	15
	Iowithla River	95	16
¹ Eastern Bristol Bay	Whale Mountain Creek	95	17
	Alagnak River	95	18
North Alaska Peninsula	Frosty Creek	95	19
	Sapsuk - Nelson River	95	20
South Alaska Peninsula Kodiak	Portage Creek	95	21
	Rough Creek	95	22
Southcentral Alaska	Little Susitna River weir	95	23
	Beartrap Creek	95	24
Southeast Alaska	Chilkat River - 24Mile	95	25
	North Arm Creek	95	26
British Columbia	Kitimat River	95	27
	Kitwanga River	95	28
Washington	Nisqually River Hatchery	95	29
	Elwha River	95	30

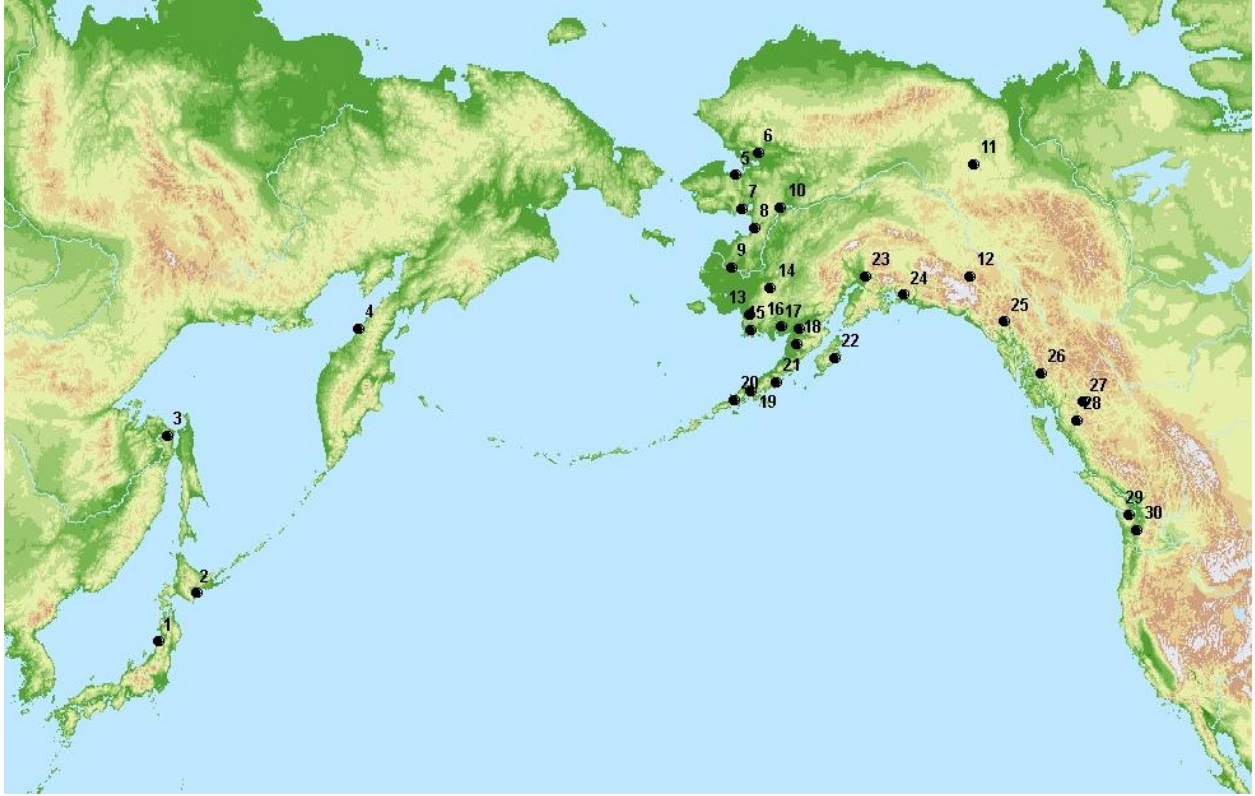
324

325 ¹ Populations in the Coastal Western Alaska (CWAK) Region.

326

327

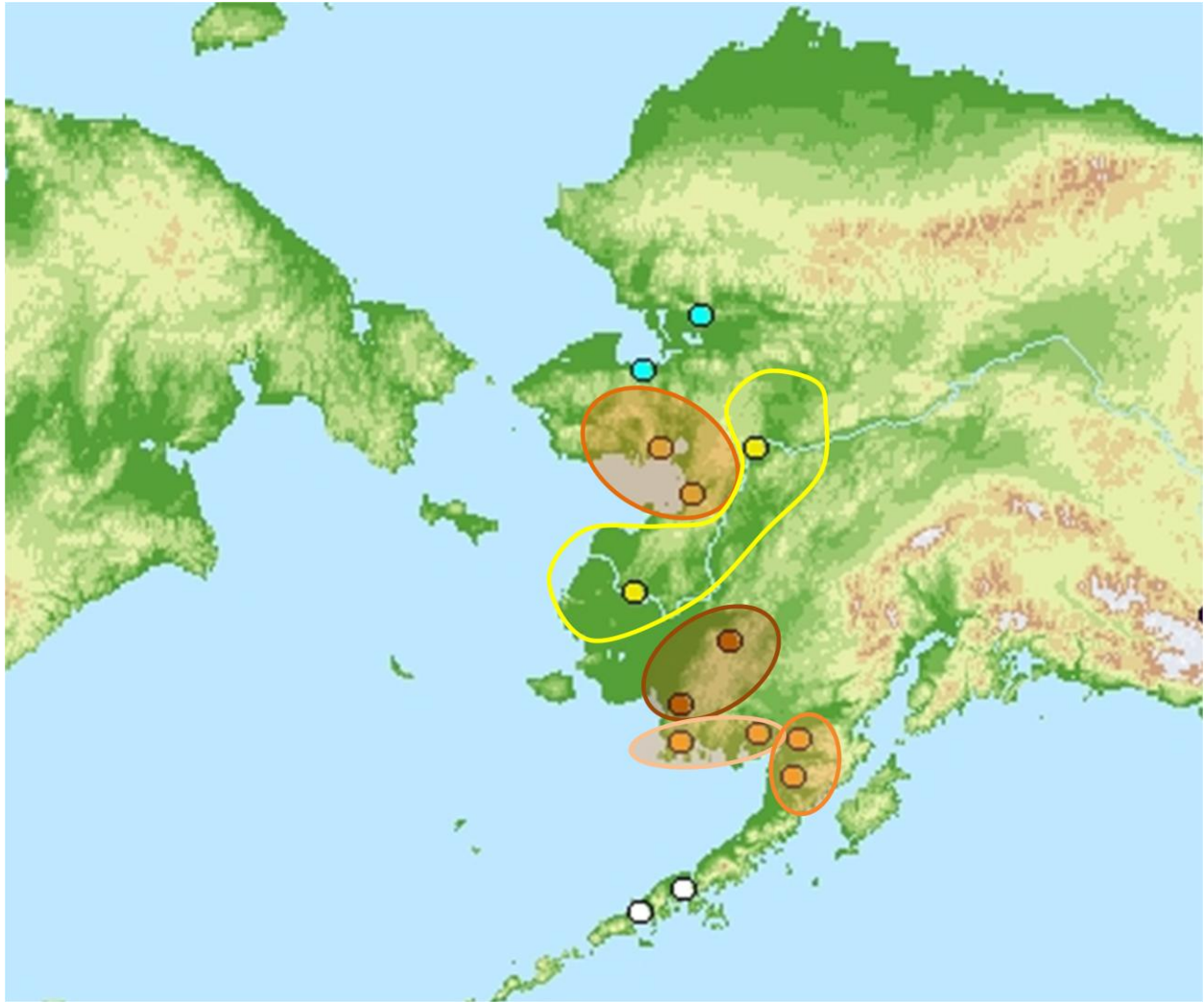
328



329
330

331 Figure 1. Map of chum salmon populations used in SNP selection process.

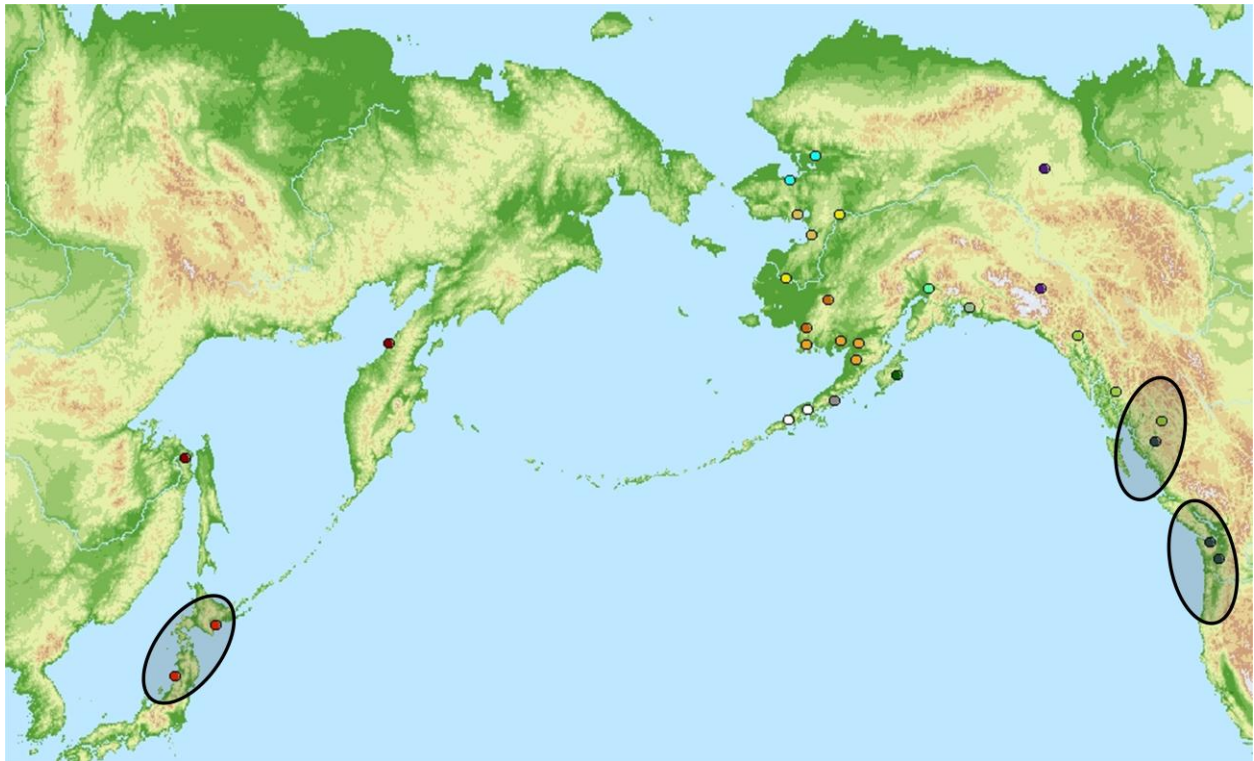
332
333



334

335

336 Figure 2. Locations of chum salmon collections within western Alaska. The five regions within
337 Western Coastal Alaska to be measured using overall F_{ST} are indicated by the ellipses.



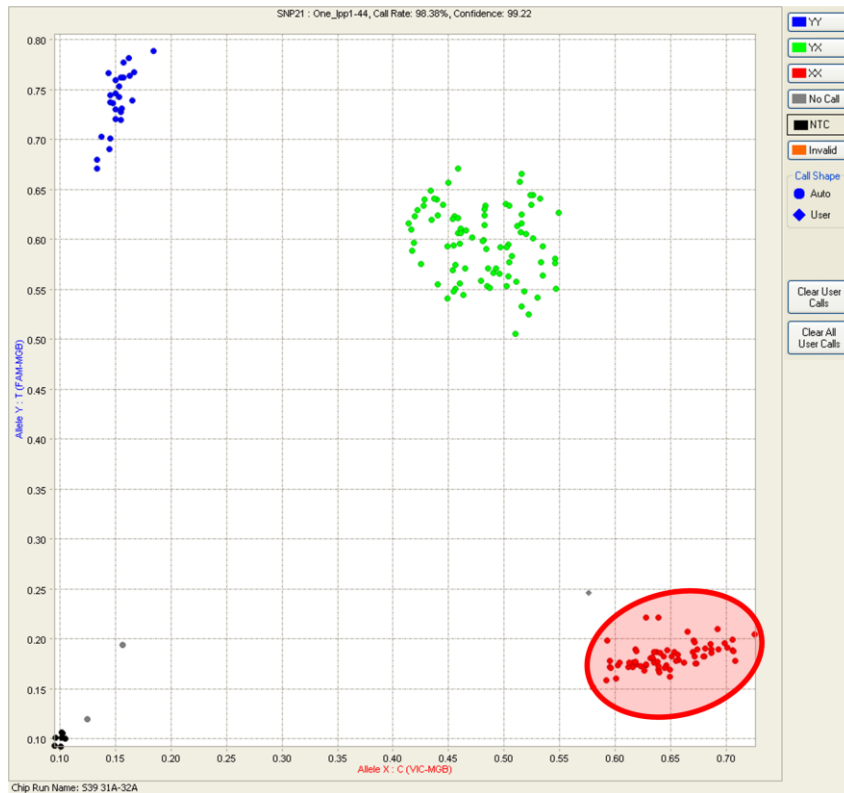
338

339

340 Figure 3. Chum salmon populations used in SNP selection process highlighting the three
341 population pairs (in ovals) of chum salmon chosen to measure F_{ST} within regions of interest to
342 research groups outside of Alaska.

343

344



345

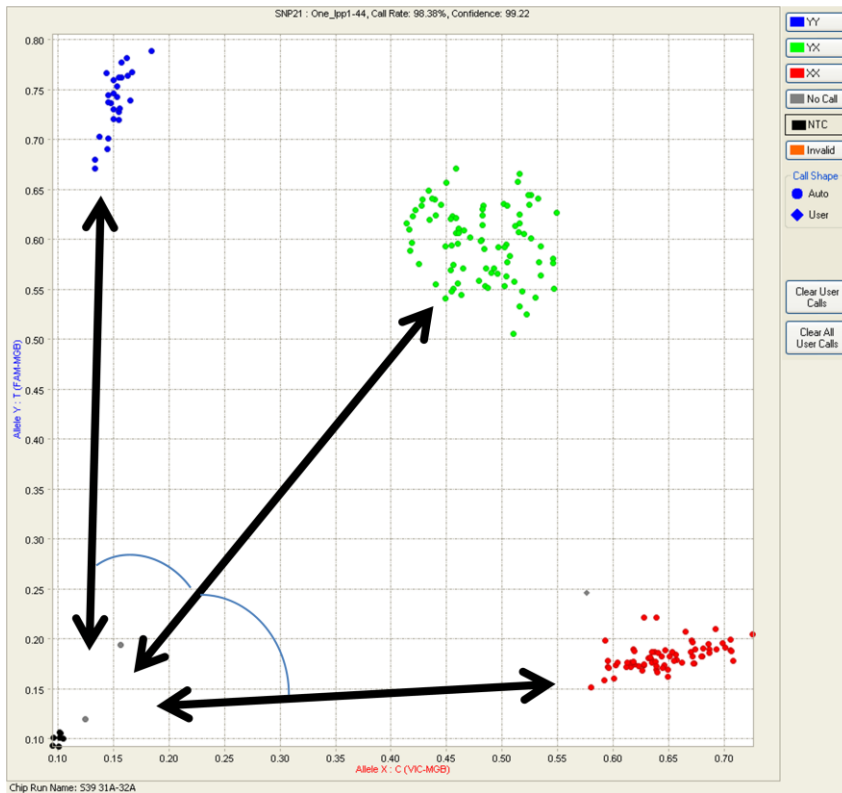
346

347 Figure 4. Screen capture of a scatter plot from genotyping software. Each point represents a
 348 single fish. The three clusters represent each possible genotype (TT homozygote - blue, TC
 349 heterozygote - green, and CC homozygote - red). The size of the shaded area for the CC
 350 homozygote distribution is an indication of cluster tightness.

351

352

353

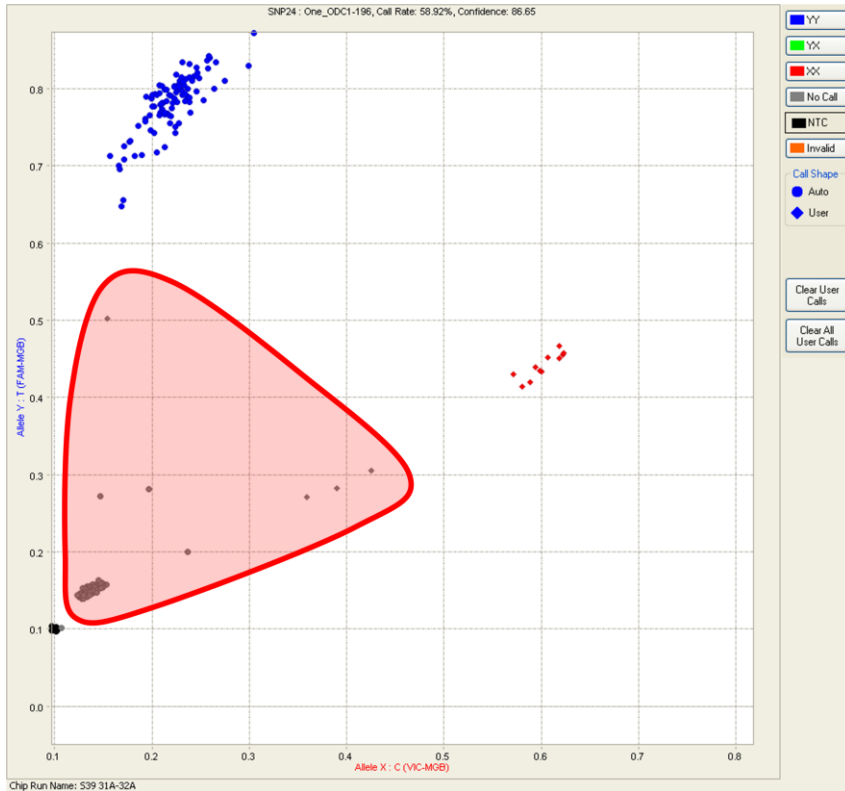


354

355

356 Figure 5. Screen capture of a scatter plot from genotyping software. Each point represents a
 357 single fish. The three clusters represent each possible genotype (TT homozygote - blue, TC
 358 heterozygote - green, and CC homozygote - red). The angle between the double-ended arrows is
 359 an indication of cluster alignment.

360



361

362

363 Figure 6. Screen capture of a scatter plot from genotyping software. Each point represents a
 364 single fish. The three clusters represent each possible genotype (TT homozygote - blue, TC
 365 heterozygote - green, and CC homozygote - red). The red shaded area represents fish for which
 366 the assay failed.

367