

Title: Parentage SNP selection – SEAK chum
Authors: K. Shedd, T. H. Dann, C. Habicht, and W. D. Templin
Date: May 27, 2014

Version: 1.0

1 **Abstract**

2 Uncertainty about the impact of hatchery salmon on the productivity and sustainability of natural
3 stocks in Prince William Sound (PWS) and Southeast Alaska (SEAK) was the impetus for the
4 Alaska Hatchery Research Program (AHRP). One major portion of this project is designed to
5 use genetic data to perform parentage analysis in order to create pedigrees and assess the impact
6 on fitness (productivity) of natural pink *Oncorhynchus gorbuscha* and chum *O. keta* salmon
7 stocks due to straying of hatchery pink and chum salmon. **Single nucleotide polymorphisms**
8 **(SNPs)** have been identified as the marker type for the parentage analysis. Markers are being
9 developed for pink salmon and markers are available for chum salmon. However, the marker
10 suite for chum salmon has yet to be determined. Here we describe our intended process to select
11 the set of SNPs for chum salmon that provides the maximum resolution possible for parentage
12 analysis to meet the objectives of the AHRP.

13 **Background of AHRG**

14 Extensive ocean-ranching salmon aquaculture is practiced in Alaska by private non-profit
15 corporations (PNP) to enhance common property fisheries. Most of the approximately 1.7B
16 juvenile salmon PNP hatcheries release annually are pink salmon in Prince William Sound
17 (PWS) and chum salmon in Southeast Alaska (SEAK; Vercesi 2013). The large scale of these
18 hatchery programs has raised concerns among some that hatchery fish may have a detrimental
19 impact on the productivity and sustainability of natural stocks. Others maintain that the potential
20 for positive effects exists. ADF&G convened a Science Panel (Alaska Hatchery Research
21 Group; AHRG) whose members have broad experience in salmon enhancement, management,
22 and natural and hatchery fish interactions. The AHRG was tasked with answering three priority
23 questions:

- 24 I. *What is the genetic stock structure of pink and chum salmon in each region (PWS and*
25 *SEAK)?*

¹ This document serves as a record of communication between the Alaska Department of Fish and Game Commercial Fisheries Division and other members of the Alaska Hatchery Research Group. 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.

- 26 II. *What is the extent and annual variability in straying of hatchery pink salmon in PWS and*
27 *chum salmon in PWS and SEAK?*
- 28 III. *What is the impact on fitness (productivity) of natural pink and chum salmon stocks due*
29 *to straying of hatchery pink and chum salmon?*

30 **Introduction**

31 *Measuring the Impact on Fitness*

32 To answer the third question, we need to know the origin and pedigree of each fish captured in
33 select streams across multiple generations. **Origin** refers to the type of early life-history habitat
34 (hatchery or natural) that a fish experienced. **Pedigree** refers to the family relationship among
35 parents and offspring. ‘**Ancestral origin**’ refers to the origin of an individual’s ancestors (e.g.,
36 two parents of a single origin [hatchery/hatchery or natural/natural] or two parents of mixed
37 origin [hatchery/natural]). These ancestral origins can be determined by combining information
38 from three sources: identification of hatchery origin from otolith marks, pedigree from genetic
39 data, and age from scales (for chum salmon from SEAK). By pairing these data within fish and
40 across generations, we can estimate **reproductive success (RS)** among cross types (i.e. hatchery-
41 hatchery, hatchery-natural, and natural-natural origin crosses). The AHRG is using the **relative**
42 **reproductive success (RRS)** of hatchery-origin fish to natural-origin fish as the measure of
43 *fitness in this study* (Shedd et al. 2014).

44 *Problem: Which Markers to Use for Parentage Analysis*

45 The reconstruction of pedigrees via parentage analysis using genetic markers is based on simple
46 **Mendelian inheritance**, where an offspring inherits one of two alleles from each parent. While
47 the concept is simple, the implementation of exclusion-based parentage analysis can be
48 challenging in open systems, where not all parents are sampled, and genetic information is
49 limited and/or subject to genotyping error. For this purpose, there are a wide variety of statistical
50 likelihood methods that utilize either a frequentist-likelihood or Bayesian approach to assess the
51 probability of parent-offspring relationships. Nevertheless, any parentage analysis is subject to
52 the limitations of the genetic marker set employed. Thus, it is important to 1) select informative
53 markers; 2) select robust markers that produce highly accurate and consistent genotypes under a
54 wide range of tissue qualities, and 3) determine the requisite number of markers for successful
55 parentage analysis.

56 While **microsatellites** have historically been the marker-type of choice for parentage analysis due
57 to their high variability and general availability, SNPs have recently received increased attention
58 due to their high-throughput screening, low genotyping error rates, and transferability among
59 laboratories. With current technology at the ADF&G Gene Conservation Laboratory (GCL),
60 genotyping cost per locus for microsatellites is an order of magnitude higher than for SNPs.
61 Theoretical work has shown that a set of 60-100 SNPs with minor allele frequency (MAF) ≥ 0.3
62 allows for accurate pedigree reconstruction of large populations that contain thousands of
63 potential mothers, fathers, and offspring (Anderson and Garza 2006). This theoretical work has

64 been confirmed by empirical studies that have compared parentage analysis with both
65 microsatellites and SNPs (Hauser et al. 2011, Tokarska et al. 2009). Hauser et al. (2011)
66 compared 11 highly variable microsatellites specifically chosen for parentage analysis to 80
67 SNPs originally designed for genetic stock identification (GSI; high among-population
68 variation). Over half of the SNPs had a MAF < 0.2, a level below which SNPs rapidly lose
69 power in parentage analysis (Anderson and Garza 2006). Despite the limitations of the SNP
70 marker set used by Hauser et al. (2011) with respect to parentage analysis, the authors found that
71 assignment success was always higher for SNPs than for microsatellites across different
72 parentage analysis software programs.

73 The GCL has 188 SNPs available for chum salmon (Table 1). These 188 SNPs have been
74 previously narrowed down to 96 SNPs (the maximum number of SNPs that can be run on a
75 single high-throughput SNP chip), however this set of 96 SNPs was optimized for GSI in
76 western Alaska using high-quality samples as part of The Western Alaska Salmon Stock
77 Identification Program (WASSIP) (DeCovich et al. 2012), not parentage analysis in Southeast
78 Alaska using carcass tissues. In order to make the final selection of the best 96 SNPs for
79 parentage analysis for the AHRP, the GCL proposes to determine the performance for all 188
80 SNPs on a sample from all 4 pedigree streams and then empirically determine the set of SNPs
81 required for optimal success in parental assignments of 2014 alevin to 2013 adults in Fish Creek.

82 *Goals of Technical Document*

83 Two goals of this technical document are to:

- 84 1) Propose and document the method for selecting markers to be used in parentage analysis.
 - 85 i. Determine population genetic summary statistics for all 188 SNPs for the 4 chum
86 salmon pedigree streams sampled in SEAK.
 - 87 ii. Determine laboratory performance for all 188 SNPs for chum salmon carcasses
88 sampled in SEAK.
 - 89 iii. Determine the required number of SNPs necessary for robust, accurate parentage
90 analysis of alevin and adult chum salmon in SEAK using Fish Creek as the test
91 population.
- 92 2) Request a decision by the AHRG on these methods prior to August 2014.

93 **Methods**

94 *Phase 1: Ranking SNPs*

95 **Suitability of SNPs for parentage analysis for AHRP: All 188 SNP markers will be**
96 **assayed in 95 randomly selected adult individuals sampled in 2013 from each of the**
97 **4 pedigree streams (Figures**

98 I. Figure 1).

- 99 1. Unranked measures: Measures in this section will be given veto power and markers
 100 will be discarded if they do not pass the following tests.
- 101 a. Hardy-Weinberg Equilibrium (HWE): Conformance to HWE will be
 102 measured with Genepop version 4.0.11 (Rousset 2008). Markers out of HWE
 103 at $\alpha = 0.05$ in any of the 4 populations or exhibiting overall significance,
 104 measured across all 4 populations, at $\alpha = 0.01$ will be dropped. An overall p-
 105 value will be calculated according to Fisher's method for combined
 106 probability test.
 - 107 b. Linkage Disequilibrium: Linkage disequilibrium will be measured with
 108 Genepop version 4.0.11 (Rousset 2008). Marker pairs that exhibit linkage
 109 disequilibrium at $\alpha = 0.05$ in 3 or more of populations examined will be
 110 considered "associated" and the SNP with the lesser average MAF among
 111 populations will be dropped.
 - 112 c. Laboratory Performance: Only markers that have an overall relative scoring
 113 success rate of >80% (relative to the best-performing marker, to account for
 114 poor tissue quality) and a discrepancy rate of <2% across all 4 populations
 115 will be retained.
- 116 2. Ranked measures: The measures in this section of the selection process will be
 117 scored between 0 and 1 (worst to best) using the equation:

$$118 \quad \text{score} = \frac{2 \times (\text{mean MAF})}{(1 + \text{SD of MAF})} \quad (\text{Eqn. 1})$$

- 119 a. "Mean MAF" is the mean MAF calculated across the 4 pedigree streams and
 120 is our primary metric of interest,
- 121 b. "SD of MAF" is the standard deviation of MAF among the 4 pedigree
 122 streams. This attributes a "cost" to including markers that are highly variable
 123 among populations (i.e., useful in some but not others).
- 124 c. Each SNP that passed the unranked measures above will be assigned a rank
 125 based upon its score for this measure. This order of ranks will be used in
 126 subsequent measures below.

127 *Phase 2: Evaluating SNP Sets*

- 128 II. Empirical test of SNPs for parentage analysis: The SNP markers that pass the unranked, veto
 129 portion of Phase 1 will then be assayed for all remaining adults collected from Fish Creek in
 130 2013 and all alevin collected from Fish Creek in 2014 in order to perform parentage analysis,
 131 as these are currently the only chum samples available for parentage analysis.
- 132 1. Parentage analysis: Parentage analysis will be performed by assigning alevin to
 133 adults using the highest ranked markers (according to Equation 1) in sets of SNPs that
 134 are efficiently screened using GCL laboratory methods: 24, 48, 96, 120, 144, and all
 135 SNPs marker sets.

- 136 2. Cost/benefit analysis: We will examine the relationship between increasing number
137 of markers used in a set (cost) and success in parentage analysis. Dependent
138 variables will include:
- 139 a. Number of offspring assigned to two parents
 - 140 b. Number of offspring assigned to one parent
 - 141 c. Number of parents with one or more successfully assigned offspring
 - 142 d. Mean log-odds (LOD) score ratios between most likely parent and the next
143 most likely parent from Cervus3 (Hauser et al. 2011) as a measure of
144 confidence in parentage assignment.
 - 145 i. Cervus3, a likelihood parentage assignment program, provides LOD
146 score for each parent-offspring assignment.
 - 147 ii. A LOD score is the natural log of the likelihood ratio (i.e. probability
148 of a putative-parent-offspring pair being related divided by the
149 probability that they are unrelated).
 - 150 iii. For a given offspring, the LOD score is computed for multiple putative
151 parents.
 - 152 iv. The ratio of LOD scores between the most likely parent and the second
153 most likely parent gives a proxy for the level of confidence in
154 assigning the most likely parent.
 - 155 v. The distribution of LOD ratios can be compared between marker sets
156 to assess the level of confidence in correct parent assignments.
 - 157 e. Parentage error rates (putative)
 - 158 i. Error rates defined by comparing assignments of a set to assignments
159 made with all available markers (Gold Standard).
 - 160 ii. Type I error – assigning an untrue parent
 - 161 iii. Type II error – failing to assign a true parent when the true parent is
162 present in the sample

163 *Final Considerations*

- 164 III. Final considerations: The candidate SNPs will be ordered from best to worst with respect to
165 the measures in the ranked portion of Phase 1 (Equation 1). Given that the GCL is optimized
166 to use 96 SNP markers in a set, the top 96 candidates from Phase 1 will be selected, unless
167 Phase 2 suggests that equal assignment power can be obtained with 48 SNPs or that more
168 than 96 SNPs are necessary to acquire adequate power.

169 **Questions for the AHRG**

- 170 1. Are the proposed methods for ranking markers appropriate and sufficient? Are there
171 other considerations that should be assessed as well?
- 172 2. Are the proposed methods for determining marker set appropriate?

173

AHRG Review and Comments

174 *This technical document was discussed at the December 12, 2014 meeting of the AHRG. In*
175 *addition it was reviewed by email exchange prior to the meeting.*

176 The proposed methods are acceptable.

177 This document is acceptable to the AHRG.

178

References

- 179 Anderson, E. C., and J. C. Garza. 2006. The Power of Single-Nucleotide Polymorphisms for Large-Scale Parentage
180 Inference. *Genetics* 172:2567-2582. <http://www.genetics.org/cgi/content/abstract/172/4/2567>
- 181 DeCovich, N. A., J. R. Jasper, S. M. Turner, C. Habicht, and W. D. Templin. 2012. Western Alaska Salmon Stock
182 Identification Program Technical Document 23: Chum salmon SNP selection results. Alaska Department of
183 Fish and Game, Division of Commercial Fisheries, Regional Information Report 5J12-25, Anchorage.
184 <http://www.adfg.alaska.gov/FedAidPDFs/RIR.5J.2012.25.pdf>
- 185 Elfstrom, C. M., C. T. Smith, and L. W. Seeb. 2007. Thirty-eight single nucleotide polymorphism markers for high-
186 throughput genotyping of chum salmon. *Molecular Ecology Notes* 7(6):1211-1215 (5).
187 <http://www3.interscience.wiley.com/journal/120808786/abstract?CRETRY=1&SRETRY=0>
- 188 Hauser, L., M. C. Baird, R. Hilborn, L. S. Seeb, and J. E. Seeb. 2011. An empirical comparison of SNPs and
189 microsatellites for parentage and kinship assignment in a wild sockeye salmon (*Oncorhynchus nerka*)
190 population. *Molecular Ecology Resources* 11(Supplement 1):13.
191 <http://www.ncbi.nlm.nih.gov/pubmed/21429171>
- 192 Petrou, E. L., L. Hauser, R. S. Waples, W. D. Templin, D. Gomez-Uchida, and L. W. Seeb. 2013. Secondary contact
193 and changes in coastal habitat availability influence the nonequilibrium population structure of a salmonid
194 (*Oncorhynchus keta*). *Molecular Ecology* 22(23):5848-5860 (13).
195 <http://onlinelibrary.wiley.com/doi/10.1111/mec.12543/pdf>
- 196 Rousset, F. 2008. GENEPOP 007: a complete re-implementation of the GENEPOP software for Windows and
197 Linux. *Molecular Ecology Resources* 8(1):103-106 (4).
198 [http://www.arlis.org:2074/ehost/pdfviewer/pdfviewer?vid=4&hid=10&sid=c9d59b7a-94df-4e14-925f-
200 df0bb3d382e3%40sessionmgr4](http://www.arlis.org:2074/ehost/pdfviewer/pdfviewer?vid=4&hid=10&sid=c9d59b7a-94df-4e14-925f-
199 df0bb3d382e3%40sessionmgr4)
- 201 <http://genepop.curtin.edu.au/>
- 202 Seeb, J. E., C. E. Pascal, E. D. Grau, L. W. Seeb, W. D. Templin, T. Harkins, and S. B. Roberts. 2011. Transcriptome
203 sequencing and high-resolution melt analysis advance single nucleotide polymorphism discovery in
204 duplicated salmonids. *Molecular Ecology Resources* 11(2):335-348 (14).
205 <http://onlinelibrary.wiley.com/doi/10.1111/j.1755-0998.2010.02936.x/abstract>
- 206 Shedd, K. R., T. H. Dann, C. Habicht, and W. D. Templin. 2014. Alaska Hatchery Reserach Program Technical
207 Document 1: Defining relative reproductive success: which fish count? ADF & G Technical Document:10.
- 208 Smith, C. T., J. Baker, L. Park, L. W. Seeb, C. M. Elfstrom, S. Abe, and J. E. Seeb. 2005a. Characterization of 13
209 single nucleotide polymorphism markers for chum salmon. *Mol. Ecol. Notes*:259-262.
210 http://www.researchgate.net/publication/227610256_Characterization_of_13_single_nucleotide_polymorphism_markers_for_chum_salmon/file/79e4150ed8f04e96d3.pdf
- 211 Smith, C. T., C. M. Elfstrom, J. E. Seeb, and L. W. Seeb. 2005b. Use of sequence data from rainbow trout and
212 Atlantic salmon for SNP detection in Pacific salmon. *Molecular Ecology* 14:4193-4203.
213 http://doc.nprb.org/web/publication/project_0205-0303_seeb_mol_ecol_2005.pdf
- 214 Tokarska, M., T. Marshall, R. Kowalczyk, J. Wójcik, C. Pertoldi, T. Kristensen, V. Loeschcke, V. R. Gregersen, and C.
215 Bendixen. 2009. Effectiveness of microsatellite and SNP markers for parentage and identity analysis in
216 species with low genetic diversity: the case of European bison. *Heredity* 103(4):326-332.
217 <http://www.nature.com/hdy/journal/v103/n4/full/hdy200973a.html>

218 Vercessi, L. 2013. Alaska salmon fisheries enhancement program 2012 annual report. Alaska Department of Fish
219 and Game, Fishery Management Report No. 13-05, Anchorage.
220 <http://www.adfg.alaska.gov/FedAidPDFs/FMR13-05.pdf>

221

223 Table 1.–188 chum salmon SNPs available for use in SEAK chum salmon parentage analyses

Assay	Source ^a	Assay	Source ^a	Assay	Source ^a
Oke_PPA2-635	A	Oke_gdh1-234	B	Oke_ras1-249	A
Oke_ACOT-100	B	Oke_gdh1-62	B	Oke_RFC2-618	C
Oke_AhR1-278	A	Oke_GHII-3129	A	Oke_RH1op-245	C
Oke_AhR1-78	A	Oke_glr1-78	B	Oke_ROA1-209	B
Oke_APOB-60	B	Oke_GNMT-100	B	Oke_RPN1-80	B
Oke_arf-319	C	Oke_GnRH-373	E	Oke_RS27-81	B
Oke_ATP5L-105	B	Oke_GPDH-191	C	Oke_RS27-94	B
Oke_ATP5L-248	B	Oke_GPH-105	A	Oke_RS9-379	B
Oke_azin1-90	B	Oke_GPH-78	A	Oke_RSPRY1-106	D
Oke_brd2-118	D	Oke_H2AX-72	B	Oke_serpin-140	C
Oke_brp16-65	B	Oke_hmgb1-66	B	Oke_slc1a3a-86	B
Oke_CATB-60	B	Oke_hnRNPL-239	A	Oke_sylc-90	B
Oke_ccd16-77	D	Oke_HP-182	A	Oke_TCP1-78	A
Oke_CCT3-143	A	Oke_HSP90BA-299	A	Oke_TCTA-202	B
Oke_CCT3-220	A	Oke_IGFI.1	C	Oke_TCTA-99	B
Oke_CD123-62	B	Oke_il-1racp-67	C	Oke_Tf-278	A
Oke_CD81-108	B	Oke_IL8r2-406	B	Oke_thic-84	B
Oke_CD81-173	B	Oke_IL8r-272	E	Oke_txnrd1-74	B
Oke_cjo57-86	B	Oke_KPNA2-87	A	Oke_u0602-244	D
Oke_CKS1-70	B	Oke_lactb2-71	B	Oke_U1001-79	D
Oke_CKS1-94	B	Oke_lamp2-138	B	Oke_U1002-165	D
Oke_CKS-389	E	Oke_LAMP2-186	B	Oke_U1002-262	D
Oke_CO1A1-72	B	Oke_mcf2-86	B	Oke_U1008-83	D
Oke_CO1A1-76	B	Oke_METK2-97	B	Oke_U1010-154	D
Oke_col1a2-62	B	Oke_mgll-49	B	Oke_U1010-251	D
Oke_Cr30	E	Oke_MLRN-63	B	Oke_U1012-241	D
Oke_Cr386	E	Oke_Moesin-160	C	Oke_U1012-60	D
Oke_ctgf-105	A	Oke_nc2b-148	B	Oke_U1015-255	D
Oke_CTR2-82	B	Oke_ND3-69	E	Oke_U1016-154	D
Oke_DBLOH-79	B	Oke_ndub3-58	B	Oke_U1017-52	D
Oke_DCXR-87	B	Oke_NHERF-123	B	Oke_U1018-50	D
Oke_DM20-548	E	Oke_NHERF-54	B	Oke_U1019-218	D
Oke_e2ig5-50	B	Oke_NUPR1-70	B	Oke_U1020-75	D
Oke_EF2-394	B	Oke_PDIA3-475	B	Oke_U1021-102	D
Oke{EIF4EB	C	Oke_PDIA3-82	B	Oke_U1022-114	D
Oke_eif4g1-43	B	Oke_pgap-111	B	Oke_U1022-139	D
Oke_f5-71	B	Oke_pgap-92	B	Oke_U1023-147	D
Oke_FANK1-166	B	Oke_pnrc2-78	B	Oke_U1024-113	D
Oke_FANK1-96	B	Oke_psm9-188	B	Oke_U1025-135	D
Oke_FBXL5-61	B	Oke_psm9-57	B	Oke_U1027-89	D
Oke_gdh1-191	B	Oke_rab5a-117	B	Oke_U1028-100	D

Assay	Source ^a	Assay	Source ^a	Assay	Source ^a
Oke_U1031-132	D	Oke_U2026-64	B	Oke_U2057-80	B
Oke_U1103-150	B	Oke_U2029-79	B	Oke_U212-87	C
Oke_u1-519	E	Oke_U2031-37	B	Oke_U216	C
Oke_U2001-629	B	Oke_U2032-74	B	Oke_u217-172	C
Oke_U2002-200	B	Oke_U2033-122	B	Oke_u200-385	C
Oke_U2003-142	B	Oke_U2034-55	B	Oke_U302-195	A
Oke_U2005-62	B	Oke_U2035-54	B	Oke_U502-241	A
Oke_U2006-109	B	Oke_U2037-76	B	Oke_U503-272	A
Oke_U2007-190	B	Oke_U2038-32	B	Oke_U504-228	A
Oke_U2010-94	B	Oke_U2040-77	B	Oke_U505-112	A
Oke_U2011-107	B	Oke_U2041-84	B	Oke_U506-110	A
Oke_U2015-151	B	Oke_U2042-61	B	Oke_U507-286	A
Oke_U2016-118	B	Oke_U2043-51	B	Oke_U507-87	A
Oke_U2017-87	B	Oke_U2045-43	B	Oke_U509-219	A
Oke_U2019-112	B	Oke_U2047-49	B	Oke_U510-204	A
Oke_U202	C	Oke_U2048-91	B	Oke_U511-271	A
Oke_U2020-51	B	Oke_U2049-99	B	Oke_U514-150	A
Oke_U2021-86	B	Oke_U2050-101	B	Oke_UBA3-245	D
Oke_U2022-101	B	Oke_U2052-56	B	Oke_uqcrfs-69	B
Oke_U2023-99	B	Oke_U2053-60	B	Oke_XBP1-82	B
Oke_U2024-93	B	Oke_U2054-58	B	Oke_zn593-152	B
Oke_U2025-86	B	Oke_U2056-90	B		

225 a A= Elfstrom et al. 2007; B= Petrou et al. 2013; C= Smith et al. 2005b; D= Seeb et al. 2011; and E= Smith et al.
 226 2005a.
 227

Figures

229 Figure 1.—Map of 4 chum salmon pedigree streams in SEAK.

