

Title: Experimental design of pink salmon SNP development
Authors: T. H. Dann, C. Habicht, W. D. Templin, L. W. Seeb and J. E. Seeb
Date:

Version: 1.0

Abstract

The fitness component of the Alaska Hatchery Research Program requires single nucleotide polymorphisms (SNPs) that are currently unavailable. The Gene Conservation Laboratory and the Seeb Laboratory of the University of Washington are using restriction site associated DNA sequencing to develop SNPs for pink salmon. Here we describe the experimental design we are using to develop SNPs, provide preliminary sequencing results, and propose SNP selection criteria to the Science Panel. DNA from 665 individual pink salmon sampled from 17 populations has been sequenced. Average retained reads was 1.6M for the 190 individuals that have had both rounds of sequencing, suggesting that we will achieve adequate sequence coverage to accurately estimate allele frequencies and identify variable SNPs useful for the fitness study. We propose a series of gating criteria and ranking measures to select SNPs for the fitness study and seek feedback from the Science Panel.

Background of AHRP

Extensive ocean-ranching salmon aquaculture is practiced in Alaska by private non-profit corporations (PNP) to enhance common property fisheries. Most of the approximately 1.7B juvenile salmon that PNP hatcheries release annually are pink salmon in Prince William Sound (PWS) and chum salmon in Southeast Alaska (SEAK; Vercessi 2014). The large scale of these hatchery programs has raised concerns among some that hatchery fish may have a detrimental impact on the productivity and sustainability of natural stocks. Others maintain that the potential for positive effects exists. To address these concerns ADF&G convened a Science Panel for the Alaska Hatchery Research Program (AHRP) whose members have broad experience in salmon enhancement, management, and natural and hatchery fish interactions. The AHRP was tasked with answering three priority questions:

- I. *What is the genetic stock structure of pink and chum salmon in each region (PWS and SEAK)?*
- II. *What is the extent and annual variability in straying of hatchery pink salmon in PWS and chum salmon in PWS and 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 Science Panel of the Alaska Hatchery Research Program. 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

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 AHRP 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* (Tech Doc 1 – Shedd et al. 2014). The current design is for RRS to be
44 estimated in six populations: Erb Creek, Hogan Bay, Paddy Creek, Spring Creek, Gilmour Creek
45 and Stockdale Creek.

46 *Identification of genetic markers*

47 Estimating relative reproductive success via parentage analyses on the scale of this research
48 program necessitates single nucleotide polymorphisms (SNPs). While microsatellites have
49 historically been the marker-type of choice for parentage analysis due to their high variability
50 and general availability, SNPs have recently received increased attention due to the potential for
51 high-throughput screening, low genotyping error rates, and transferability among laboratories.
52 With current technology at the ADF&G Gene Conservation Laboratory (GCL), genotyping cost
53 per locus for microsatellites is an order of magnitude higher than for SNPs. With all parents
54 sampled, theoretical work has shown that a set of 60-100 SNPs with minor allele frequency
55 (MAF) > 0.3 allows for accurate pedigree reconstruction of large populations that contain
56 thousands of potential mothers, fathers, and offspring (Anderson and Garza 2006). This
57 theoretical work has been confirmed by empirical studies that have compared parentage analysis
58 with both microsatellites and SNPs (Hauser et al. 2011, Tokarska et al. 2009, Anderson 2012).
59 Hauser et al. (2011) compared 11 highly variable microsatellites specifically chosen for
60 parentage analysis to 80 SNPs originally designed for genetic stock identification (GSI; high
61 among-population variation). Over half of the SNPs had a MAF < 0.2, a level below which
62 SNPs rapidly lose power in parentage analysis (Anderson and Garza 2006). Despite the
63 limitations of the SNP marker set used by Hauser et al. (2011) with respect to parentage analysis,
64 the authors found that assignment success was always higher for SNPs than for microsatellites
65 across different parentage analysis software programs.

66 Accurately and confidently assigning offspring to parents requires many independently assorting
67 alleles at genetic markers that are variable within the population being studied. Recent
68 simulation work indicates that ~192 independent alleles with MAF greater than 0.3 can resolve
69 parent-offspring relationships in study conditions expected for the AHRP (large populations, not
70 all parents sampled; Shedd et al. 2015). The fitness study will also need to analyze many
71 thousands of individuals in the laboratory, so genetic markers with high throughput capabilities
72 are required. These factors make it clear that SNPs are the genetic marker to use for the fitness
73 study. Other benefits of SNPs include the ability to select from 10,000s of potential markers, the
74 reliability of genotype calls as a direct measure of sequence, reduced expense to genotype, and
75 transferability among laboratories. Taqman assays (ADF&G's current SNP genotyping
76 methodology) have been developed for 51 pink salmon SNPs that exhibit signatures of selection
77 (University of Washington, unpublished), but these likely are not sufficient to resolve parentage
78 in AHRP fitness streams. Thus, developing a panel of 192 SNPs that resolve parentage in both
79 odd- and even-year lineages is a major objective of the fitness study.

80 *Marker development*

81 Restriction site associated DNA sequencing (RAD sequencing) has emerged as the best approach
82 for SNP discovery and is the method used to develop SNPs for AHRP (Andrews et al. 2016).
83 The approach consists of: 1) digesting DNA with restriction enzymes to isolate RAD tag sites
84 among all individuals, 2) ligating a short sequence of individual-specific bases to the cut site to
85 barcode each fragment to individual, 3) shearing DNA to reduce the length of sequence
86 fragments for sequencing, 4) generating many copies of RAD tag sites via PCR amplification, 5)
87 pooling individuals into libraries and 6) sequencing libraries to generate sufficient copies of
88 RAD tag sites to confidently genotype individuals and discover variant bases (SNPs; Baird et al.
89 2008). Benefits of this approach include the standardized method that has been well vetted and
90 documented in primary literature, the potential to select from among 10,000s of SNPs, and the
91 transferability of information from SNPs developed via RAD sequencing across studies. The
92 method has been successfully used in many taxa to address a variety of questions including
93 identification of SNPs to distinguish closely related populations of Chinook salmon and genomic
94 regions of divergence between ecotypes of threespine stickleback (Hohenlohe et al. 2010, Larson
95 et al. 2014). The method has also proven successful in identifying SNPs exhibiting parallel
96 selection between the two lineages of pink salmon (Seeb et al. 2014).

97 *Selecting markers in context of genome organization*

98 Constructing marker panels in the context of a linkage map is important to selecting independent
99 markers tailored to their application. Linkage maps describe the order and relative spacing of
100 markers along chromosomes (or linkage groups). The order and relative spacing are estimated
101 from the frequencies of recombinations between loci and are often generated by examining
102 segregation patterns of haploid or diploid families (i.e., Limborg et al. 2015). Selecting markers
103 located on different chromosomes or distant enough from one another within chromosomes
104 ensures statistical independence among markers, an assumption made by many common

105 analyses. For the purposes of parentage analyses, statistical independence of genotypes is
106 important for statistical power as tightly linked markers provide redundant information for
107 resolving parent-offspring relationships.

108 *Goals of Technical Document*

109 The goals of this technical document are to:

- 110 1) Describe the experimental design used to develop SNPs for pink salmon in Prince
111 William Sound;
- 112 2) Provide results to date and an expected timeline for the remainder of SNP development
113 process; and
- 114 3) Propose and, ask for input from the Science Panel on, the selection criteria for SNPs to be
115 used for parentage analyses.

116 **Methods**

117 *Identification of samples for RAD sequencing*

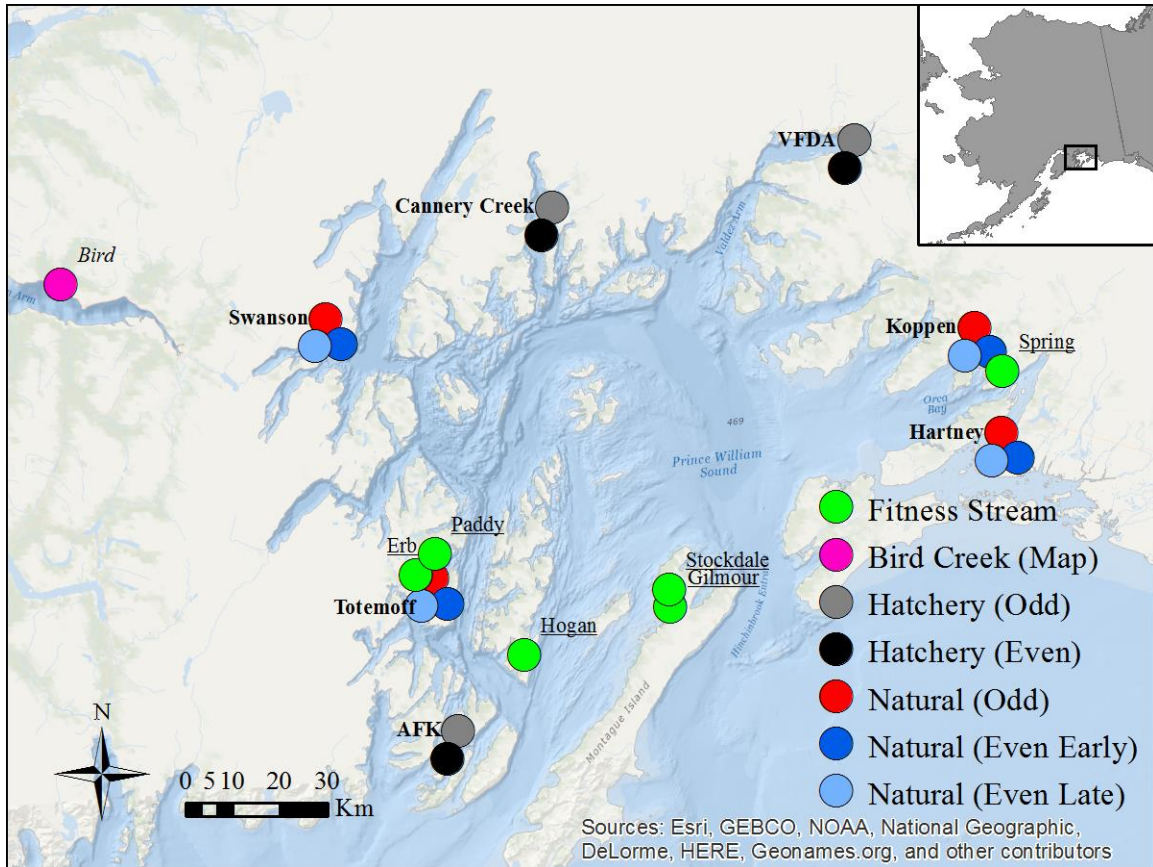
118 A primary objective of our study is to identify a panel of genetic markers that maximizes
119 parentage power in the six streams under analysis by the AHRP. Marker development for AHRP
120 would ideally use individuals from fitness streams to identify highly variable markers. However,
121 samples from individuals from nearby streams are a better choice for this particular study. The
122 samples from fitness streams are of poor quality. Tissues were sampled from only post-
123 spawning mortalities; such tissues often have degraded DNA that fail to provide quality RAD
124 data (Graham et al. 2015). As a result we identified three criteria of representative populations to
125 help select from available samples: 1) availability of quality tissue samples with paired data to
126 identify natural-origin individuals and provide the potential to identify sex-linked SNPs
127 discovered via RAD sequencing; 2) expectation of presence in fitness streams (i.e. one of the
128 three hatchery populations potentially present in fitness streams; Habicht et al. 2000); and 3)
129 expectation that population allele frequencies will be similar to natural fitness-stream
130 populations. We considered tissue samples to be of good quality if they originated from
131 spawning individuals of natural-origin pink salmon collected within the recent five years. When
132 available, we also attempted to include early and late components of pink salmon runs to
133 incorporate any temporal variability that may exist within the six fitness populations.

134 **Results**

135 *Experimental design of RAD sequencing*

136 The experimental design of our RAD sequencing effort was determined by the objectives and
137 criteria above. The GCL has 250 collections of pink salmon populations (including hatcheries)
138 from Prince William Sound; 134 of these originate from allozyme collections from the 1990s
139 that likely include hatchery fish and/or degraded DNA. Seventy-six collections were sampled in
140 2013 and 2014, 20 of which are either hatchery or fitness stream populations. Of the remaining

141 56 collections, we selected 17 for RAD sequencing based upon the objectives and criteria above.
 142 We included the three hatchery populations expected to be observed in fitness streams: Armin F.
 143 Koernig Hatchery (AFK), Cannery Creek Hatchery, and Solomon Gulch Hatchery (aka VFDA).
 144 We attempted to select collections from West and East Prince William Sound to parallel the
 145 distribution of pink salmon populations throughout the Sound and represent allele frequencies
 146 expected in fitness streams (Figure 1; Table 1).



147
 148 Figure 1. Locations of Alaska Hatchery Research Program fitness study streams (green), the population
 149 used for linkage map construction (Bird Creek from Cook Inlet), and hatchery (gray/black) and natural
 150 (red/blue) populations selected for RAD sequencing. Fitness study stream names are underlined, the
 151 name of the linkage map stream is italicized and RAD sequence stream names are in bold.

152 Table 1. Name, Gene Conservation Laboratory (GCL) code, collection date, latitude, longitude
 153 for collections used to develop SNPs. Also included are the number of samples available (Total)
 154 and included in RAD sequencing.

Collection	GCL Code	Collection Date	Latitude	Longitude	Sample size	
					Total	Sequenced
AFK Hatchery	PAFK13	8/22/2013	60.051	-148.065	200	37
Cannery Creek Hatchery	PCANN13	8/25/2013	61.019	-147.514	197	37
VFDA Hatchery	PVFDA13	8/9/2013	61.084	-146.304	200	37

Totemoff Creek	PTOTM13	8/27/2013	60.343	-148.088	96	37
Swanson Creek	PSWAN13	8/12/2013	60.849	-148.412	121	37
Koppen Creek	PKOPP13	7/30/2013	60.706	-145.898	216	37
Hartney Creek	PHART13	8/3/2013	60.502	-145.842	271	37
AFK Hatchery	PAFK14	8/25/2014	60.065	-148.065	200	37
Cannery Creek Hatchery	PCANN14	8/26/2014	61.019	-147.514	200	37
VFDA Hatchery	PVFDA14	7/31/2014	61.131	-146.348	200	37
Totemoff Creek Early	PTOTM14E	7/29/2014	60.343	-148.088	96	37
Totemoff Creek Late	PTOTM14L	8/26/2014	60.343	-148.088	102	37
Swanson Creek Early	PSWAN14E	7/28/2014	60.849	-148.412	120	37
Swanson Creek Late	PSWAN14L	8/25/2014	60.849	-148.412	125	37
Koppen Creek Early	PKOPP14E	8/2/2014	60.706	-145.898	120	35
Koppen Creek Late	PKOPP14L	8/31/2014	60.706	-145.898	124	38
Hartney Creek Early	PHART14E	8/4/2014	60.502	-145.842	46	37
Hartney Creek Late	PHART14L	8/21/2014	60.502	-145.842	109	37

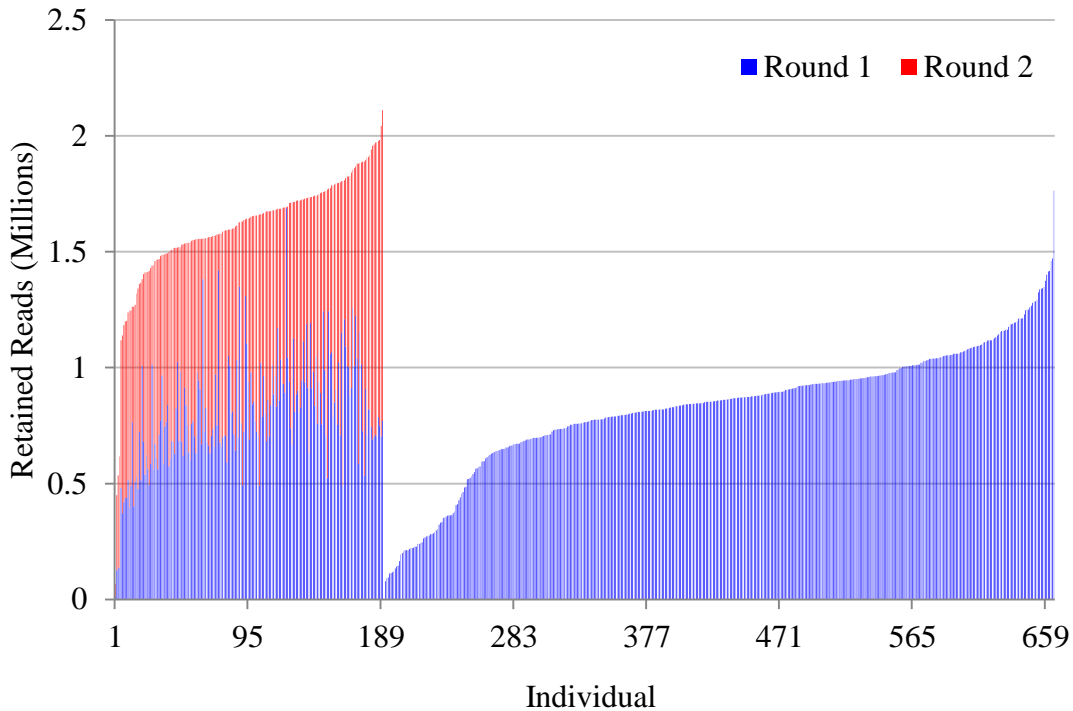
155

156

Sequencing results

157 We have completed 64% of the sequencing and expect to complete the remainder in January,
 158 2016. Our approach to collecting RAD data involves sequencing each library of 95 individuals
 159 twice at University of Oregon’s Genomics and Cell Characterization Core Facility
 160 (<https://gc3f.uoregon.edu/>). The first round of sequencing used constant volumes of extracted
 161 DNA solution from each fish. The second round of sequencing is conducted after DNA
 162 concentrations are titrated to equalize the sequence coverage (number of retained reads) among
 163 individuals based on results from the first round. Equal coverage among individuals results in
 164 better coverage among markers and subsequently better overall estimates of allele frequencies
 165 for each population. Two libraries have been sequenced twice while the other five are in the
 166 process of being re-sequenced.

167 Initial results indicate that the goal of 1.5M retained reads will be reached for each individual.
 168 Individuals from the two libraries that have had both rounds of sequencing averaged 1.6M
 169 retained reads, while individuals from the remaining 5 libraries averaged 825K retained reads
 170 after one round (Figure 2). As expected, average retained reads/collection were unequal after the
 171 first round of sequencing (e.g., PVFDA13=232K, PAFK14=1,063K) but became much more
 172 equal after the second round of sequencing for the 6 completed collections (range 1,534-1,755K;
 173 Table 2).



174

175 Figure 2. The number of retained reads for 665 pink salmon sequenced to identify SNPs for the
 176 Alaska Hatchery Research Program. Of the 665 fish, 190 have had both round 1 (blue) and 2
 177 (red) of sequencing while the remaining 475 are expected to be completed in January.

178 Table 2. Number of individuals sequenced and average retained reads after the first, second, and
 179 both rounds of sequencing.

Collection	n	Round 1	Round 2	Both
PAFK13	37	911,239	622,276	1,533,515
PCANN13	37	881,855	750,057	1,631,912
PHART13	37	848,992	793,622	1,642,613
PSWAN13	37	761,967	849,031	1,610,998
PKOPP13	37	600,574	1,008,059	1,608,633
PTOTM13	37	603,598	1,084,180	1,754,826
PVFDA13	37	232,305	NA	NA
PAFK14	37	1,063,014	NA	NA
PCANN14	37	930,839	NA	NA
PVFDA14	37	917,992	NA	NA
PHART14E	37	854,268	NA	NA
PHART14L	37	795,811	NA	NA
PKOPP14E	35	872,938	NA	NA
PKOPP14L	38	929,163	NA	NA
PSWAN14E	37	895,058	NA	NA

PSWAN14L	37	875,550	NA	NA
PTOTM14E	37	786,726	NA	NA
PTOTM14L	37	941,559	NA	NA

180

181

Proposed criteria to select SNPs for parentage purposes

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

We propose a series of gating criteria and ranking measures to select the set of 192 SNPs from the 10,000s available for the parentage panel. The gating criteria will reduce the pool of SNPs through exclusion based on adverse characteristics. As the first gating criterion, we will exclude loci that are duplicated as a result of the whole genome duplication common to the ancestor of all salmonids (Ohno 1970) to ensure that only unambiguous genotypes are used to determine parent-offspring relationships. Similarly, we will exclude loci that exhibit a minor allele frequency < 0.05 in any collection and those that fail to conform to Hardy-Weinberg expectations. Finally, to ensure independence of SNPs included in our panel, we will attempt to include only those SNPs with a minimum separation of 15 centiMorgans (cM). The cM is a measure of genetic distance between loci on a chromosome that is defined as the number of chromosomal crossovers (recombinations) expected among 100 individuals in a single generation and thus is a good measure to identify independence among loci. For example, 15 cM equates to 15 crossover events occurring for every hundred offspring in a generation between two loci. The current linkage map available for pink salmon, based upon odd-year populations from Washington State and Bird Creek, Alaska, characterizes 26 chromosomes with an average length of 129 cM (Limborg et al. 2014; University of Washington, unpublished). Given our goal of developing a panel of ~200 SNPs, we expect to select ~8 SNPs from each chromosome which translates into 16 cM spacing between SNPs if we assume similar chromosome lengths.

200

201

202

203

To choose from the remaining SNPs, we propose to use the same ranking measure proposed for the selection of chum salmon SNPs in Technical Document 2 (Shedd et al. 2014). Each SNP will be assigned a score based upon the mean and standard deviation (SD) of minor allele frequency (MAF) across the RAD collections following:

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

204

205

206

This measure standardizes scores between 0 and 1 and is an intuitive measure to base parentage power as MAF is the most important factor in a marker's power in parentage analyses (Anderson and Garza 2006).

207

208

209

210

We expect to identify hundreds of SNPs that can resolve parent-offspring relationships for AHRP and will select the final set based upon the gating and ranking criteria above as well as each marker's ability to be accurately genotyped using amplicon technology. The GCL is transitioning to the GT-seq methods (Campbell et al. 2014); ensuring that the final panel of

211 markers produces consistent sequence coverage among markers will result in more accurate
212 genotype calls.

213 **Questions for the AHRP**

214 1. Are the proposed methods for selecting SNPs for parentage analyses adequate? Do you
215 suggest other approaches to selecting SNPs?

216 **AHRP Review and Comments**

217 *This technical document has been reviewed.*

218 This document is acceptable to the AHRG.

219 There was once comment by Alex Wertheimer who stated:

220 *The approach the authors have outlined for selecting SNPs for use in the parentage analysis is*
221 *very well-justified and scientifically sound. The authors are taking advantage of advance*
222 *techniques in identifying SNPs in the genome of pink salmon, and have devised criteria for*
223 *selecting the "best" 192 SNP markers for use in the parentage analyses for the many thousands*
224 *that are available. The authors' comprehensive description of the rationale and methodology*
225 *continues the excellent job of documenting the scientific rigor brought to bear to address the*
226 *priority questions and objectives of the Alaska Hatchery Research Program. I certainly have no*
227 *recommendations on alternate approaches.*

228 **References**

229 Anderson, E.C. 2012. Large-scale Parentage Inference with SNPs: an Efficient Algorithm for Statistical
230 Confidence of Parent Pair Allocations. *Statistical applications in genetics and molecular biology* **11**(5).

231 Anderson, E.C., and Garza, J.C. 2006. The power of single-nucleotide polymorphisms for large-scale
232 parentage inference. *Genetics* **172**(4): 2567-2582.

233 Andrews, K.R., Good, J.M., Miller, M.R., Luikart, G., and Hohenlohe, P.A. 2016. Harnessing the power
234 of RADseq for ecological and evolutionary genomics. *Nat Rev Genet* **advance online publication**.

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

238 Campbell, N.R., Harmon, S.A., and Narum, S.R. 2014. Genotyping-in-Thousands by sequencing (GT-
239 seq): A cost effective SNP genotyping method based on custom amplicon sequencing. *Molecular Ecology*
240 *Resources*: n/a-n/a.

241 Graham, C.F., Glenn, T.C., McArthur, A.G., Boreham, D.R., Kieran, T., Lance, S., Manzon, R.G.,
242 Martino, J.A., Pierson, T., Rogers, S.M., Wilson, J.Y., and Somers, C.M. 2015. Impacts of degraded
243 DNA on restriction enzyme associated DNA sequencing (RADSeq). *Molecular Ecology Resources* **15**(6):
244 1304-1315.

- 245 Habicht, C., Simpson, E.M., and Seeb, J.E. 2000. Broodstock acquisition and release sites for hatcheries
246 producing pink salmon in Prince William Sound. Alaska Department of Fish and Game, Regional
247 Information Report No. 5J00-07.
- 248 Hauser, L., Baird, M., Hilborn, R., Seeb, L.W., and Seeb, J.E. 2011. An empirical comparison of SNPs
249 and microsatellites for parentage and kinship assignment in a wild sockeye salmon (*Oncorhynchus nerka*)
250 population. *Molecular Ecology Resources* **11**: 150-161.
- 251 Hohenlohe, P.A., Bassham, S., Etter, P.D., Stiffler, N., Johnson, E.A., and Cresko, W.A. 2010. Population
252 genomics of parallel adaptation in threespine stickleback using sequenced RAD Tags. *PLOS Genetics* **6**(2
253): e1000862.
- 254 Larson, W.A., Seeb, J.E., Pascal, C.E., Templin, W.D., and Seeb, L.W. 2014. Single-nucleotide
255 polymorphisms (SNPs) identified through genotyping-by-sequencing improve genetic stock identification
256 of Chinook salmon (*Oncorhynchus tshawytscha*) from western Alaska. *Canadian Journal of Fisheries and*
257 *Aquatic Sciences* **71**(5): 698-708.
- 258 Limborg, M.T., Waples, R.K., Allendorf, F.W., and Seeb, J.E. 2015. Linkage Mapping Reveals Strong
259 Chiasma Interference in Sockeye Salmon: Implications for Interpreting Genomic Data. *G3:
260 Genes|Genomes|Genetics* **5**(11): 2463-2473.
- 261 Limborg, M.T., Waples, R.K., Seeb, J.E., and Seeb, L.W. 2014. Temporally Isolated Lineages of Pink
262 Salmon Reveal Unique Signatures of Selection on Distinct Pools of Standing Genetic Variation. *Journal*
263 *of Heredity* **105**(6): 741-751.
- 264 Ohno, S. 1970. Enormous diversity in genome sizes of fish as a reflection of nature's extensive
265 experiments with gene duplication. *Transactions of the American Fisheries Society* **99**(1): 120-&.
- 266 Seeb, L.W., Waples, R.K., Limborg, M.T., Warheit, K.I., Pascal, C.E., and Seeb, J.E. 2014. Parallel
267 signatures of selection in temporally isolated lineages of pink salmon. *Mol Ecol* **23**(10): 2473-2485.
- 268 Shedd, K.R., Dann, T.H., Habicht, C., and Templin, W.D. 2014. Alaska Hatchery Research Program
269 Technical Document 1: Defining relative reproductive success: which fish count? ADF & G Technical
270 Document.
- 271 Shedd, K.R., Dann, T.H., Habicht, C., and Templin, W.D. 2014. Alaska Hatchery Research Program
272 Technical Document 2: Parentage SNP selection - SEAK chum. ADF & G Technical Document.
- 273 Tokarska, M., Marshall, T., Kowalczyk, R., Wojcik, J.M., Pertoldi, C., Kristensen, T.N., Loeschcke, V.,
274 Gregersen, V.R., and Bendixen, C. 2009. Effectiveness of microsatellite and SNP markers for parentage
275 and identity analysis in species with low genetic diversity: the case of European bison. *Heredity* **103**(4):
276 326-332.
- 277 Vercesi, L. 2014. Alaska salmon fisheries enhancement program 2013 annual report. Fishery
278 Management Report. Alaska Department of Fish and Game, Anchorage.
- 279