# Dense SNP panels resolve closely related Chinook salmon populations

SCHOLARONE™
Manuscripts

1

2

3     Dense SNP panels resolve closely related Chinook salmon populations

4

5     Garrett J. McKinney[1a*], Carita E. Pascal[1b], William D. Templin[2c], Sara E. Gilk-Baumer[2d], Tyler H. Dann[2e],

6     Lisa W. Seeb[1f], James E. Seeb[1g]

7     [1] School of Aquatic and Fishery Sciences, University of Washington, 1122 NE Boat Street, Box 355020,

8     Seattle WA 98195-5020, USA.

9     [2] Alaska Department of Fish and Game, 333 Raspberry Road, Anchorage, AK, USA  99518

10    *Corresponding author: Garrett J. McKinney (email: gjmckinn@uw.edu, phone: 1-765-430-3272)

11    .Emails: [a] gjmckinn@uw.edu; [b]cpascal@uw.edu; [c] bill.templin@alaska.gov; [d] sara.gilk@alaska.gov;

12    [e]tyler.dann@alaska.gov; [f]lseeb@uw.edu; [g]jseeb@uw.edu

13

14    **Running head:** SNP markers for Chinook salmon

15    **Key words:** Chinook salmon, genetic stock identification, amplicon sequencing, RADseq, GT-seq,

16

17

18

## Abstract

Chinook salmon are migratory fish that are highly valued for subsistence, sport, and commercial fisheries throughout their native range. Populations of Chinook salmon in Western Alaska have exhibited long-term declines, leading to restrictions on harvests. Management priorities require greater resolution for genetic stock identification (GSI) than is available with current methods. We leveraged RADseq, TaqMan, and GT-seq data originating from multiple sources, collected through time, to develop a set of GT-seq panels containing 1,092 SNPs that improved GSI resolution in Western Alaska for at-sea and in-river sampling. We generated a dense linkage map with to ensure that markers selected for panels spanned the entire genome. In addition, we identified multiple RADseq markers that were associated with sex; these aligned to a 5cM region on the sex chromosome. Finally, we developed a bioinformatic pipeline to streamline analysis of GT-seq data that is capable of genotyping microhaplotypes and paralogs, both of which can improve GSI resolution over traditional single-SNP data. Our panels and pipeline provide tools for management agencies to rapidly and easily analyze large-scale genotyping projects.

**Introduction**

Genomic data have become a central feature of management and conservation of fish populations (e.g.,

Bernatchez et al. 2017, Sylvester et al. 2017). Conservation applications enabled by genome-wide data sets

include studies of adaptation, genotype-by-environment interactions, inbreeding and outbreeding depression,

or loss of adaptive variation (Allendorf et al. 2010). Genomic data enable traceability of escapees from fish

farms (Pritchard et al. 2016, Holman et al. 2017), important because such escapees provide a major threat to

genetic variability and sustainability of wild populations (Bolstad et al. 2017, Forseth et al. 2017). Finally,

effective harvest management of migratory species such as salmon requires accurate identification of unique

populations as they mix along their migratory corridors (Dann et al. 2013, Meek et al. 2016), and SNP

panels identified from genome-wide data sets frequently provide management and conservation solutions

(McKinney et al. 2017, Beacham et al. 2018).

It is now relatively easy to generate thousands of markers to address conservation and management questions

either through whole-genome or reduced representation sequencing; however, high-throughput genotyping of

thousands of samples is still prohibitive in terms of cost and time. One solution is to distill these datasets to

subsets of informative loci that can be genotyped using amplicon sequencing. This allows a rapid and cost-

effective method for genotyping thousands of individuals for hundreds of loci.

Chinook salmon are migratory fish that are important for ceremonial, subsistence, sport, and commercial

fisheries throughout their native range in Pacific Rim drainages from the Kamchatka Peninsula, in the

western Pacific Ocean, to Central California, USA, in the eastern Pacific Ocean (Healey 1991). The species

is culturally significant among indigenous tribes; some celebrate the first Chinook salmon caught each year

57   with spiritual ceremonies.  Their large size (up to 50 kg) and fighting ability make Chinook salmon a prized

58   sportfish, and their size and flesh quality make them a highly valued commercial and sport fish.  These

59   factors have collided to complicate management and allocation of migrating Chinook salmon among user

60   groups (Miller 1993, Lin et al. 1996, Gisclair 2009) and among nations (deReynier 1998, Walsh 1998).

61

62   Many populations of Chinook salmon have exhibited long-term declines throughout their native range (e.g.,

63   Schoen et al. 2017, Siegel et al. 2017).  Population declines in the Eastern Bering Sea have led to restrictions

64   on both subsistence and commercial harvest, and individuals are becoming smaller and younger at maturity

65   (Ohlberger et al. 2018).  These conservation challenges are compounded by the fact that many fisheries

66   harvest multiple populations (i.e., mixed stock fisheries) of Chinook salmon of differing abundances.

67   Effective management requires the ability to identify the components of individual stocks harvested in

68   mixtures so that less productive stocks can be protected and harvest can target concentrations of productive

69   stocks (Beacham et al. 2008).  In addition to management of mixed stock fisheries, a main driver for

70   distinguishing stocks of Western Alaska Chinook salmon is the desire to characterize the composition of

71   bycatch in the walleye pollock fishery (Templin et al. 2011); up to 60% of the Chinook salmon bycatch can

72   originate from Western Alaska Chinook salmon (Myers and Rogers 1988, Myers et al. 2009).

73

74   Populations of Chinook salmon throughout the coastal areas of the eastern Bering Sea show little genetic

75   differentiation when genotyped by traditional panels of up to 96 Single Nucleotide Polymorphisms (SNPs)

76   and subdivide into only five genetically identifiable reporting groups for mixture analysis or genetic

77   assignment tests  (Larson et al. 2014a).  These reporting groups consist of Norton Sound and three reporting

78   groups related to differentiation within the Yukon River (Upper Yukon River, Middle Yukon River, Lower

79   Yukon River); the final reporting group combines all populations from the Kuskokwim Drainage/Bristol Bay

4

80  region.  Analyses that genotype 1000s of SNPs demonstrate that larger panels of SNPs show promise for

81  further resolving these populations (restriction site associated DNA sequencing or RADseq; Larson et al.

82  2014b, see also Sylvester et al. 2017).  At the same time, techniques have emerged that enable high-

83  throughput genotyping of 100s of information-rich SNPs (e.g., GT-seq; Campbell et al. 2015), suitable for

84  conservation and management applications.  Our objective was to use RADseq to discover and screen 1000s

85  of SNPs to find informative markers and then use them to construct GT-seq panels to provide further

86  resolution of Chinook salmon populations in Western Alaska.

87

88  We report the development and testing of a set of marker panels containing a total of 1,092 existing and

89  newly ascertained SNPs for discriminating stocks of Chinook salmon.  Existing SNPs, already tested in other

90  applications, originated from three sources: a 299-SNP GT-seq panel adapted for use with Columbia River,

91  USA, populations (Hess et al. 2016) and developed in part from SNPs ascertained in Western Alaska (Smith

92  et al. 2005a, Smith et al. 2005b, Smith et al. 2005c); a 96 SNP TaqMan panel developed for genetic stock

93  identification in Western Alaska (Larson et al. 2014a); and 178 RADseq SNPs developed for population

94  discrimination in Cook Inlet, Alaska (Dann et al. 2018).  Newly ascertained SNPs originated from a dense

95  RADseq screen of individuals from populations spanning the Kuskokwim and Nushagak river drainages in

96  Western Alaska.  The final set of 1,092 SNPs increased resolution for identifying components of potential

97  harvest mixtures, providing a total of eight reporting groups for Chinook salmon in Western Alaska where

98  only five existed before.  In addition, simulations of sampling local fisheries supported the discrimination of

99  three reporting groups each within the Kuskokwim and Nushagak drainages.

100

101  Finally, we developed a computational pipeline for processing and genotyping of GT-seq data that is able to

102  genotype multiple-SNP haplotypes and markers with different levels of ploidy.  Data handling is increasingly

103 challenging with increasing numbers of SNPs; early interest suggests that the computational pipeline will

104 expedite processing of GT-seq, RAPTURE (Ali et al. 2016), and similar data sets.

105

106 Our work demonstrates how an increasingly large number of SNPs can be collated and evaluated for high

107 throughput analyses; we also present recommendations for further streamlining the discovery and evaluation

108 phases that should be widely applicable to species of management and conservation interest. These new

109 SNP panels will aid in management of Alaska Chinook salmon by improving the ability of managers to

110 discriminate among stocks of interest during harvest. Lessons learned in this study will be widely applicable

111 to other studies where genetic stock identification is applied to mixtures of migrating populations.

112

113 **Materials and Methods**

114 *RADseq data*

115 We conducted SNP discovery on RADseq data from 13 populations (Figure 1, Figure 2A-C; Figure 3; Table

116 1). Raw data for five of the populations were available from previous studies (Larson et al. 2014b (2

117 populations), McKinney et al. 2018 (3 populations)) and downloaded from Dryad (doi:10.5061/dryad.rs4v1)

118 and NCBI (SRA SRP129894) (Table 1). Additional RAD sequencing was conducted on 48 fish per

119 population from 8 new populations (Table 1). In total, nine populations from the Kuskokwim drainage and

120 four populations from Bristol Bay (Togiak River and Nushagak drainage) were used for SNP discovery.

121 Full materials and methods for RAD sequencing are available in supplemental file S1.

122 We upgraded the McKinney et al. (2016) genetic map for Chinook salmon to provide a framework to help

123 ensure that markers chosen for our panels were distributed across the genome. Genetic maps improve the

124    power to identify genes and gene regions important in population differentiation and adaptation.   The

125    McKinney et al. (2016) map was a consensus of maps derived from populations from Puget Sound,

126    Washington, and likely missed variation present in northern stocks (see Templin et al. 2011).   We mapped

127    additional variation detected in five haploid families that were created from females sampled from Ship

128    Creek (Anchorage, Alaska) to produce a new consensus map (supplemental file S2).

129    *Sex Locus*

130    We attempted to incorporate loci that could be used to determine the sex of immature fish into our SNP

131    panels. Knowledge of sex of individual fish and sex ratios in migrating cohorts provides useful demographic

132    information, but identification of reliable DNA sex markers has been elusive (Von Bargen et al. 2015).  Sex

133    data were only available from 23 fish from the Togiak River, so putative sex-associated loci were instead

134    identified by examining population-level genotype patterns.  In salmon, males are the heterogametic sex and

135    females are the homogametic sex  (Thorgaard 1977), and in Chinook salmon the sex-determining region is

136    on chromosome 17 (Ots17, Phillips et al. 2013).  Assuming equal numbers of males and females in a

137    population, a sex locus should be revealed by the presence of high heterozygosity (~50%) but only two

138    genotypic classes, one heterozygous and one homozygous, with a resulting minor allele frequency of 0.25.

139    We identified putative sex-associated loci as those with heterozygosity between 45% and 55%, minor allele

140    frequency between 0.20 and 0.30.  Up to 5% of individuals with the alternate homozygous genotype were

141    allowed to account for genotyping error.  Putative sex loci were placed on the linkage map (below) to

142    determine if they co-locate with a previously identified sex-associated locus from Chinook salmon

143    originating from the Marblemount Hatchery, Skagit River drainage, in Washington State (University of

144    Washington, unpublished data).

145    *GT-seq panel construction*

146 The populations examined in this study collapsed into a single reporting group with marker panels currently

147 used for management. We simulated mixed stock analyses from the RADseq dataset using *GSIsim*

148 (Anderson et al. 2008, Anderson 2010) to determine the number of markers necessary to subdivide the

149 existing Kuskokwim/Bristol Bay reporting group into major river drainages. A training holdout leave-one-

150 out (THL) approach was used to minimize bias in assignment accuracy (Anderson 2010). Samples within a

151 population were randomly assigned to training and holdout datasets. The training set was used to choose

152 high $F_{ST}$ loci for panel testing; the holdout set was used to evaluate panel accuracy against the baseline

153 (training plus holdout samples). Simulations were done in intervals of 250 markers, up to 1,500 markers,

154 after which they were done in intervals of 1,000. Simulation results showed that approximately 1,000

155 markers were necessary to resolve reporting groups into major river drainages. GT-seq panel development

156 targeted a final set of 1,000 SNPs for population analyses.

157 We chose to partition SNPs into four GT-seq panels of about 300 SNPs each. Experience has shown that up

158 to 300 amplicons is a workable number to optimize panel performance (Beacham et al. 2018, McKinney et

159 al. 2018). Also, compartmentalizing SNPs allowed us to potentially capture existing data genotyped on other

160 platforms and shape modular panels for downstream applications on subsets of SNPs.

161 A panel of 299 existing SNPs (Panel 1) was available from Idaho Fish and Game that was originally

162 developed by the Columbia River Inter-Tribal Fish Commission (CRITFC, Hess et al. 2016). This panel

163 includes the majority of the 192 loci from Warheit et al. (2013) that form the basis for data collection for

164 Pacific Salmon Commission applications from California to Southeast Alaska (e.g. Clemento et al. 2014).

165 Panels 2-4 were developed by this study. Panel 2 originated from other currently available SNPs, either in

166 use by Alaska Department of Fish and Game (ADF&G) or proposed for use by ADF&G in Cook Inlet,

167 Alaska (Dann et al. 2018) to facilitate comparison to pre-existing data (including the $F_{ST}96$ SNPs of Larson

168 et al. (2014a) that were ascertained in western Alaska). Panels 3 and 4 originated from novel RADseq loci

169    ascertained in this project to specifically improve resolution within and between the closely related

170    populations in Kuskokwim and Nushagak river drainages in western Alaska.

171    Informative markers were identified using both outlier analysis and $F_{ST}$ for panels 3 and 4 (Figure 2C).

172    Outlier loci were identified using *BayeScan* with default settings (Foll and Gaggiotti 2008) as those with Q <

173    0.05. $F_{ST}$ was estimated using Genepop (Rousset 2008). Initial testing showed many outliers with very low

174    $F_{ST}$. These loci were likely false positives; they were generally only variable in one or a few populations,

175    and their absence in other populations was likely due to sampling error. Loci were re-filtered to include only

176    loci with a minimum minor allele frequency of 0.05 in three or more populations and then reanalyzed with

177    *BayeScan*.

178     High $F_{ST}$ markers from both within and between the Kuskokwim and Nushagak river drainages were

179    considered. Excess candidate markers, to allow for dropouts, were chosen as follows (Figure 2C): 1) 500

180    markers with highest $F_{ST}$ within the Kuskokwim drainage for Panel 3 ($F_{ST}$ range: 0.051-0.005, 2) 500

181    markers with highest $F_{ST}$ within the Nushagak drainage and Togiak River for Panel 4 ($F_{ST}$ range: 0.106-

182    0.006), and 3) 500 highest $F_{ST}$ markers between the Kuskokwim and Nushagak drainages ($F_{ST}$ range: 0.620-

183    0.004) for incorporation into either Panel 3 or Panel 4. Loci were split between panels 3 or 4 based on $F_{ST}$

184    within regions; panel 3 contained loci with greater $F_{ST}$ in the Kuskokwim Bay/River while panel 4 contained

185    loci with greater $F_{ST}$ in Togiak Bay/Nushagak drainage. Outlier loci identified using *BayeScan* were

186    contained within the high $F_{ST}$ locus set.

187    These 1500 loci underwent additional filtering to meet design criteria for GT-seq analysis (Figure 2D). Loci

188    with SNPs within 16 bp of the 3' end and 20 bp of the 5' end of the RAD locus were excluded to allow room

189    for primer design. Where matches were available, paired-end contigs from Larson et al. (2014b)  were used

190    to extend the 5' end of the RAD locus for primer design. Loci with low complexity or transposable element

191    sequence were identified using RepeatMasker (Smit et al. 2013) and removed. Primers were designed using

9

192   Primer3 (You et al. 2008) with default settings for loci that passed initial filters.  Primers were then aligned

193   to all loci using bowtie2 (Langmead and Salzberg 2012) to identify and remove cases where primers may

194   amplify multiple loci.  A total of 706 loci passed filters and were retained for panel optimization; 350 loci for

195   panel 3 and 356 loci for panel 4.  In combination with panels 1 and 2, a total of 1,343 loci were passed

196   forward for the initial optimization (Figure 2D).

197   *GT-seq Genotyping*

198   For panel optimization and population genotyping we developed a bioinformatic pipeline, *GTscore*

199   (https://github.com/gjmckinney/GTscore) to score both multiple SNP haplotypes (also referred to as

200   microhaplotypes; Baetscher et al. 2018) and duplicated loci; both of these locus types have been shown to

201   increase power for resolving closely related populations (Limborg et al. 2017, McKinney et al. 2017, Waples

202   et al. 2017) and were included in the GT-seq panels.  Multi-SNP haplotypes occur when multiple SNPs are

203   in the same sequence tag, resulting in haplotypes with > two alleles.  This pipeline incorporates the *polyGen*

204   algorithm (McKinney et al. 2018) that uses a maximum likelihood method for genotyping and is capable of

205   genotyping loci with any number of alleles and any level of ploidy.

206

207   *Panel Optimization*

208   Optimization of Panel 2, Panel 3, and Panel removed loci that did not amplify properly in the PCR reaction

209   and perfected PCR performance for each of the retained primer pairs (Figure 2 E).  Panel 1 had already

210   undergone extensive optimization by Hess et al. (2016) but was included in panel optimization to keep PCR

211   conditions consistent for the final sequencing run and ensure there were no interactions with loci from other

212   panels.

213    Optimization was done with 100bp paired-end sequencing in two rounds on an Illumina MiSeq (Figure 2E).

214    DNA was extracted, and sequencing libraries prepared following the methods of Campbell et al. (2015).  The

215    first round of sequencing used 48 individuals from four populations (Kogrukluk, Koktuli, Necons, and

216    Togiak rivers).  Primer performance was evaluated in this first round (see below); however, read depth for

217    most loci was too low to allow accurate genotyping.  The sample size was reduced in the second round of

218    sequencing to 24 samples from two populations (Kogrukluk and Koktuli rivers, 12 samples each) to ensure

219    adequate read depth to evaluate genotype concordance between RADseq and GT-seq.

220    After each round of sequencing, we eliminated loci that were over-amplifiers, off-target amplifiers, or cross-

221    amplifiers.  Over-amplifying loci generate excessive sequences relative to other loci in the panel and can be

222    identified by examination of ranked number of sequence reads per locus.  Off-target amplifiers generate

223    sequences that do not match the target sequence; these sequences contain either the forward primer or the

224    reverse primer (sometimes both) but do not contain the bioinformatic probe that identifies allelic variation

225    within the target sequence.  Total read counts for each locus were used to identify over-amplifiers, and

226    counts of primer and probe alignments were used to identify off-target amplifiers.  Paired-end sequencing

227    generates sequence reads from both ends of a DNA sequence.  With GT-seq, both reads in a pair should be

228    from the same locus, and the R1 read should start with the forward primer while the R2 read should start

229    with the reverse primer.  Cross-amplification occurs when a product is amplified by the forward and reverse

230    primers of two different loci.  This could occur when multiple loci are physically close or when multiple

231    regions of the genome are genetically similar.  Cross-amplifiers are identified where the R1 and R2 reads of

232    an amplicon align to different loci.  We used a custom pipeline for identifying cross-amplifiers.  Reference

233    sequence for each locus was generated by trimming the RAD consensus sequence to contain only the

234    sequence between each primer.  For each individual, GT-seq sequences were aligned to the reference

235    sequence using *GATK* (McKenna et al. 2010) and *SAMtools* (Li et al. 2009).  The resulting alignments were

11

236  processed with custom perl scripts to quantify cross-amplification per locus.  Patterns of cross-amplification

237  were visualized using network plots in R.

238  Validation of genotyping accuracy is necessary when combining results between different technologies

239  because loci do not always genotype consistently across technologies.  After the second round of sequencing,

240  we excluded loci with discordant genotypes that could not be explained by low read depth in either the

241  RADseq or GT-seq genotyping.  Extensive cross validation between RADseq and single SNP data (assayed

242  by 5'-nuclease reaction with TaqMan chemistry) genotypes was already done by Hess et al. (2016) for loci in

243  Panel 1) and Larson et al. (2014a) for loci in Panel 2.

244

*Baseline Data Set for Performance Testing and Mixture Analyses*

246  For final testing and mixture analyses, we prepared a baseline data set that targeted a sample size of 95 in

247  each of 17 major populations, spanning the Kuskokwim drainage, Kuskokwim Bay, and Bristol Bay

248  (Nushagak drainage and Togiak River) (Table 1, Figure 3), for all of the GT-seq loci that passed filters.  We

249  took three genotyping steps to accomplish this:  (1) we added TaqMan data to the RADseq data available for

250  the 48 individuals in the original 13 populations used for RADseq discovery (to account for TaqMan-origin

251  loci present in Panel 1 and Panel 2); (2) we used the four GT-seq panels to genotype up to 49 additional

252  individuals in each of the 13 discovery populations to approach the target sample size; and (3) we used the

253  four GT-seq panels to genotype 95 individuals in each of four new populations.

254  TaqMan methods were identical to those described in (Larson et al. 2014a), and GT-seq sequencing as

255  described above was conducted on an Illumina HiSeq 4000 with 1,190 loci and 270 samples per lane.  The

256  RADseq, TaqMan, and GT-seq datasets were combined and filtered in R. Samples shared between datasets

257  allowed further cross-validation of genotypes, ensuring genotype concordance across datasets.

258    A final GT-seq filtering step, examining allele ratio plots, was conducted following genotyping of all

259    samples (Figure 2F).  A histogram of allele ratios was plotted for each SNP for visual examination.

260    Singleton loci should have up to three peaks, depending on allele frequency, that are centered at 0, 0.5, and 1

261    (Figure S1A).  Duplicate loci should have up to five peaks centered at 0, 0.25, 0.5, 0.75, and 1 (Figure S1B).

262    Diverged duplicate loci should have up to three peaks either centered on 0, 0.25, and 0.5 or on 0.5, 0.75, and

263    1 (Figure S1C).  Loci that did not display distinct peaks associated with each genotype class are likely

264    amplifying off-target sequence (Figure S1D); we attempted to recover these loci by extending the

265    bioinformatic probe to exclude off-target sequence.  Loci that could not be recovered were removed from

266    further analysis.

267    A final filtering step prior to mixed-stock analysis (MSA) was necessary due to the nature of the combined

268    data and the fact that some data types were not present or scorable in the original RADseq data from the 13

269    populations.  Loci excluded for these reasons could be included in analysis of future GT-seq only datasets.

270    Duplicate and diverged duplicate loci were removed prior to MSA evaluation because they could not be

271    reliably genotyped in the RADseq data due to inadequate read depth (McKinney et al. 2018).  Allele

272    frequencies of RADseq and GT-seq data were also compared, and loci with allele frequency discrepancies

273    were removed.  Loci were also removed if their genotype rate was less than 70% or if they were

274    monomorphic.  Following locus filtering, samples with a genotype rate less than 90% were removed prior to

275    MSA evaluation.

276

277    *Modelling Mixed Stock Analysis*

278    The potential resolution of MSA was assessed using the full RADseq dataset while the accuracy of the four

279    GT-seq panels for MSA was assessed using the combined RADseq and GT-seq sample set.  For loci that

13

280    contained multiple SNPs, haplotypes were used to further improve accuracy (McKinney et al. 2017).

281    Populations were divided into reporting groups based on genetic affinities and management objectives.  The

282    desired reporting groups included Upper Kuskokwim River, Kuskokwim Bay, Togiak Bay, and Nushagak

283    River.  These groups are highly productive, important to stakeholders, are managed separately, and have

284    responded differently to environmental conditions responsible for recent declines throughout this region.

285    Mixture analysis and individual assignment was performed in *GSIsim* to explore alternative reporting group

286    configurations including region-wide and in-region MSA.  Mixture analysis included 100% simulations to

287    evaluate correct allocations of population to reporting group of origin and accuracy and precision among

288    reporting groups.  Individual assignment was conducted at the population level, with individuals assigned to

289    a population if they met a threshold of 80% probability assignment to that population; an 80% threshold has

290    been shown to provide a balance between a low false positive assignment rate and a successful assignment

291    rate (Griffiths et al. 2013).

292

293    **Results**

294    *Sequencing and Mapping Results*

295    We retained 19,435 loci that were scored in 761 individuals to evaluate utility for population discrimination.

296    Final sample size per population ranged from 32 to 56 with 46-48 in most populations (Table 1).   Population

297    pairs showed low overall $F_{ST}$ (average 0.003) consistent with previous studies in this region (Templin et al.

298    2011).  Results for all pairwise $F_{ST}$ comparisons are listed in Table S1.

299    Three putative sex loci were identified (Table 2). Previously, we identified a sex-associated RAD locus

300    (*RAD93920*) for Chinook salmon in the Marblemount Hatchery in Washington State within the Pacific

301    Northwest of North America (University of Washington, unpublished data); this locus was located 8 cM

302 from the centromere of *Ots17* which is the sex chromosome in Chinook salmon (Phillips et al. 2013). Two

303 of the three loci identified in this study (*RAD67724, RAD29719*) co-located with the previously identified

304 sex-locus near the centromere of *Ots17* at 4.7 cM and 9.4 cM. The other putative sex locus (*RAD27492*)

305 could not be placed on the linkage map. The accuracy of the three putative sex loci varied from 78% to 91%

306 when compared to observed sex for 23 Togiak River samples (Table 2). The SNP associated with sex in the

307 Marblemount population (*RAD93920*) was invariant in all Alaska populations.

308 A total of 15,930 loci that were scored in 233 individuals were retained for linkage mapping in the five

309 families. Sample size per family ranged from 31 to 62 (Supplementary File S1). The population and linkage

310 mapping datasets require different filtration steps so there are loci that were unique to each dataset. A total

311 of 12,140 loci were common to both datasets. The Alaska linkage map was 2,874.02 cM long and contained

312 15,798 loci; 13,084 were singleton (non-duplicated) and 2,714 were duplicate (Supplemental File S2). A

313 total of 7,946 loci mapped in the Alaska families were present on the previous map of McKinney *et al.*

314 (2016) that originated from Washington State; the low degree of shared markers is likely a due to a

315 combination of both the different pools of standing genetic variation in the source populations (Templin et al.

316 2011) and the low number of families. The combined linkage map was 3,003.36 cM long and contained

317 23,715 loci; 19,762 loci were singleton and 3,953 loci were duplicate (Table 3; Supplemental File S2). A

318 total of 630 of the newly ascertained RADseq loci that passed filters for inclusion in GT-seq panels (see

319 below) were present on the linkage map. These were distributed approximately evenly across the linkage

320 groups with an average of 18.5 GT-seq markers per linkage group and a range of 4-34 markers per linkage

321 group (Supplemental File S3). Linkage groups differed in size; the proportion of markers on each linkage

322 group that were included in GT-seq panels averaged 2.1% with a range of 1.1%-5.7%.

323

324 *GT-seq Marker Selection and Optimization*

15

325    We followed several steps of marker selection leading to GT-seq panel optimization.  SNP location within

326    the RAD tag was a limiting factor.  Approximately 25% of the 1,500 RAD markers originally selected for

327    panel design had to be discarded; the requirement of 16 bp on the 3' end of the SNP or 20 bases on the 5'

328    end was a major limiting factor for primer design from the original sequence of 94bp in length.  For SNPs

329    that were within 20 bp of the 5' end of the sequence, paired-end contigs allowed primer design past the

330    original 94 base limitation.  Minor losses of SNPs to the panels were also due to transposable element

331    annotation (4%) and identification as repetitive elements (3%).

332    A total of 1,343 loci passed initial filtering criteria and were developed into GT-seq loci for panels 2-4

333    (Figure 2) including 18 of 54 loci identified as outliers by *BayeScan*.   Two rounds of test sequencing and

334    optimization followed to remove loci due to over-amplification, excessive off-target sequence or cross-

335    amplification between primers, or genotype discrepancies between RADseq and GT-seq genotyping.  A total

336    of 1,204 loci originating from the three new panels and that of Hess et al. (2016) were genotyped for all GT-

337    seq samples (Figure 2F).

338

339    *GT-seq Genotyping and Standardization of Datasets*

340    Two additional filtering steps were conducted following genotyping.  Individuals were removed from the

341    analysis if they likely originated from a different population of origin (detailed in File S1).  The genotype

342    data also allowed us to examine the allele ratio plots for fit to expected genotype distributions; a total of 112

343    loci were removed if allele ratios did not fit the expected distributions (e.g. Figure S1D).  The final set of

344    GT-seq panels contained 1,092 loci (Figure 2G).

345    Because our study comparisons required standardization among three sources of genotypes (RADseq, GT-

346    seq, TaqMan), a final quality control step to combine the datasets prior to MSA was necessary.  Loci were

16

347  excluded if they were absent in one of the datasets (84), were duplicate or diverged duplicate loci (70),

348  exhibited greater than expected allele frequency differences between RADseq and GT-seq datasets (17),

349  exhibited a low genotype rate (34) (Figure S2), or were monomorphic (40) (Table 4).  Finally, samples were

350  excluded if they exhibited a low genotype rate (Figure S3).  A total of 847 loci and 1,545 individuals were

351  available for evaluation with *GSIsim* following the standardization steps; 14% of these loci contained

352  haplotype data.

353

354  *Mixture Analysis*

355  We first evaluated the power of the full set of more than 15,000 RAD loci.  The full dataset increased

356  resolution for Chinook salmon demonstrating > 95% accuracy to five reporting groups (Figure 3A).

357  Individuals were proportionally assigned to populations and proportional assignments were summed for

358  populations within a reporting group.  The Kuskokwim and Nushagak drainages could be resolved with the

359  full RADseq dataset.  The Kuskokwim drainage could be split into three reporting groups: Upper

360  Kuskokwim, Kuskokwim River, and Kuskokwim Bay (Figure 3A).  Additional reporting groups included a

361  combined Togiak/Goodnews reporting group and a Nushagak drainage reporting group.

362  We then evaluated the subset of 847 loci in the four GT-seq panels.  We evaluated accuracy of four different

363  reporting group scenarios of varying scales based on management objectives.  Scenarios include: 1) all

364  populations with reporting groups of Upper Kuskokwim, Kuskokwim River/Bay, Togiak/Goodnews, and

365  Nushagak rivers (fine-scale bycatch scenario, Figure 3B), 2) all populations with reporting groups of Upper

366  Kuskokwim River, combined Kuskokwim/Nushagak, and combined Togiak/Goodnews rivers (broad-scale

367  bycatch scenario, Figure 3C).  In-region scenarios include: 3) Kuskokwim drainage populations with

368  reporting groups of Upper Kuskokwim River, Kuskokwim River, and Kuskokwim Bay (Kuskokwim

369     Bay/River scenario, Figure 3D), and 4) Togiak Bay/Nushagak drainage populations with reporting groups of

370     Togiak, Iowithla/Stuyahok, and Koktuli rivers (Togiak Bay/Nushagak River scenario, Figure 3E). For some

371     scenarios Togiak and Goodnews rivers were combined into a reporting group on the basis of genetic

372     similarity even though they are not part of the same drainage. Reporting groups in each scenario had >90%

373     accuracy with two exceptions, the Nushagak River reporting group in the fine-scale bycatch scenario (72%)

374     and the Kuskokwim Bay reporting group in the Kuskokwim Bay/River scenario (87%) (Figure 3). Mean

375     accuracy and 95% range for mixture estimates under each reporting group scenario are listed in Table S2 for

376     reporting groups as a whole and for populations within reporting groups.

377     *Individual Assignment*

378     Individual assignment using the full RADseq dataset showed >95% self-assignment accuracy for all

379     populations with the exception of Kwethluk and Togiak rivers (Table S3). Individual assignment using the

380     GT-seq panels revealed that upriver populations for both the Kuskokwim and Nushagak drainages had >90%

381     self-assignment accuracy with the exception of the Takotna River population (Table S4, S5, S6).

382     Populations in the lower reaches of the Kuskokwim and Nushagak drainages misassigned to populations

383     both within and between drainages while Togiak and Goodnews rivers tended to assign to themselves or to

384     each other. Pitka Fork, Tatlawiksuk, Necons, Togiak, and Koktuli river populations all exhibited >90%

385     accuracy for either at-sea or in-region sampling.

386

387     **Discussion**

388     Genotype data are now routinely a centerpiece in the mosaic of tools used by conservation practitioners for

389     population assessment and sustainability planning. Until recently, the limited number of genetic markers

390     available often limited the power of genotype data to resolve populations even though other biological

18

391    information suggested that more resolution should be possible.  The introduction of RADseq and other

392    reduced-representation sequencing (RRS) protocols greatly increased the availability of genetic markers,

393    followed in many situations by improved resolution of closely related populations after genotyping hundreds

394    or thousands of loci (see RADseq data in Larson et al. 2014b, Candy et al. 2015).   But the cost, limited

395    analysis pipelines, and relatively slow throughput of RRS data often preclude their use in management

396    situations where data from thousands of individuals (cf., Gilbey et al. 2017) or real-time data (cf., Dann et al.

397    2013) are needed.  Amplicon sequencing approaches promise an intermediate solution where hundreds of

398    loci may be rapidly and cost-effectively genotyped in thousands of individuals, offering conservation

399    alternatives unavailable until now (Campbell et al. 2015, Beacham et al. 2018).

400    To evaluate resolution of populations we opted to build a baseline data set by adding GT-seq genotypic data

401    to the RADseq data available from SNP discovery efforts from this study and from two previous studies

402    (Larson et al. 2014b, McKinney et al. 2018).  This approach appeared to be a practical use of a large amount

403    of existing data when we started; however, two of the GT-seq panels incorporated data originating from

404    TaqMan derived loci, requiring new TaqMan genotyping to backfill those same loci missing in the RADseq

405    data sets.  While robust data emerged after careful cross validation, a better and ultimately more cost-

406    efficient choice would have been to regenotype all samples with identical GT-seq panels.

407    Genetic maps have become increasingly common and can be used as a foundation to integrate genomic

408    resources for gene annotation and population genomic analyses (McKinney et al. 2016).  However, genetic

409    maps originating from one lineage or geographic location often don't include a high proportion of

410    polymorphic markers informative for a distant lineage or geographic region.  We were able to leverage

411    existing map resources for Chinook salmon from the Pacific Northwest and add variation from families that

412    originated from Alaska to create a much denser map with much improved coverage for Alaska.   As a result,

413    74% of the markers could be placed on the combined linkage map generated in this study while only 39%

414    could be placed on the previous Washington-based linkage map.

415    The dense genetic map also allowed us to investigate the location of sex-associated loci.  The sex-associated

416    loci identified in this study were located in a 5 cM window of chromosome 17 which has been previously

417    identified as the sex chromosome in Chinook salmon (Phillips et al. 2013).  These loci flank a previously

418    identified sex marker from Marblemount Hatchery Chinook salmon (University of Washington, unpublished

419    data).  The co-location of sex associated markers from multiple projects suggest that the sex determining

420    gene is near this region of chromosome 17; however, the markers identified herein showed a maximum

421    accuracy of 91% suggesting incomplete linkage between these SNPs and the true sex determining region.

422    This is consistent with results from other molecular markers developed to assign sex to immature Chinook

423    salmon which have displayed inconsistent accuracy when tested in populations throughout the species range

424    (Nagler et al. 2001, Nagler et al. 2004, Chowen and Nagler 2005, Von Bargen et al. 2015).

425

426    *GT-seq panel*

427    The GT-seq panels increased MSA accuracy for Western Alaska Chinook salmon compared to earlier

428    analyses.  We were able to split the Kuskokwim/Bristol Bay reporting group identified by Larson et al.

429    (2014a) into three reporting groups (Upper Kuskokwim, Kuskokwim/Nushagak, Togiak/Goodnews) with

430    >95% accuracy.  These panels also show utility for regional population discrimination with three reporting

431    groups for the Kuskokwim Bay/River scenario and three reporting groups for the Togiak Bay/Nushagak

432    River scenario.  Individual assignment accuracy was generally >90% for upper river populations in both the

433    Kuskokwim and Nushagak rivers.

434    A primary goal for fisheries management in the Kuskokwim River/Bristol Bay region is the discrimination of

435    populations from the Kuskokwim and Nushagak drainages.  Splitting the Kuskokwim/Nushagak reporting

436    group by drainage yielded MSA accuracy of 93% for Kuskokwim drainage populations but only 72% for

437    Nushagak drainage populations.  A potential cause of this inaccuracy is the large difference in sample size

438    between the two reporting groups.  The Kuskokwim River reporting group contained eight populations and

439    728 samples, representing most of the production from the drainage, while the Nushagak River reporting

440    group contained three populations and 252 samples. Differences in sample size can cause bias where some

441    fish from low-sample reporting groups may assign to higher-sample reporting groups (Moran and Anderson

442    2018).

443    We were also unable to utilize the full panel in testing MSA accuracy.  The panels included paralogs which

444    may increase accuracy of genetic stock identification (Gilbey et al. 2016).  We were able to successfully

445    genotype paralogs in the GT-seq portion of the dataset, and the program we used for GSI (*GSIsim*) is capable

446    of including paralogs in analysis.  However, paralogs were excluded from MSA evaluation because read

447    depth in the RADseq dataset was too low for reliable genotypes (McKinney et al. 2018).  In addition, many

448    loci from Panel 1 were excluded either due to monomorphic genotypes (34) or absence of baseline data (61).

449    Genotyping additional Nushagak River populations would help determine if the low resolution we observed

450    in that drainage is an artifact of unbalanced sampling or a true limitation of the panel.  Sequencing of

451    additional baseline samples using GT-seq will also obviate the need to include the RADseq samples in

452    analysis, allowing paralogs to be incorporated as well as the Panel 1 loci missing in the baseline.

453    It may be possible to improve the efficiency observed in this study by working to design fewer panels with

454    higher resolution loci.  The primary source of high-$F_{ST}$ marker loss was the use of SNP position thresholds at

455    both ends of the marker sequence to allow space for primer design.  Two Chinook salmon genome

456    assemblies have now been deposited in NCBI (accessions: GCA_002872995.1, GCA_002831465.1); these

21

457   can be used to design primers for these high-$F_{ST}$ loci that were otherwise lost due to SNP position.  Also, the

458   additional power gained from haplotype loci, even when $F_{ST}$ is reduced (McKinney et al. 2017, Baetscher et

459   al. 2018), was recognized after our marker selection had concluded.  The current panels include only 114 loci

460   with haplotypes.  Locus selection that enriches for haplotypes will allow smaller panels that achieve the same

461   resolution as the current four panels and may yield further increases in resolution.

462

463   *Recommendations*

464   We leveraged RADseq, TaqMan, and GT-seq data originating from multiple sources, collected through time,

465   to develop a comprehensive baseline of population data.  While ultimately successful, this involved

466   considerable effort to standardize datasets and ~20% of the loci were discarded from analysis due to

467   incompatibilities among datasets.   Regenotyping all populations using a consistent laboratory method would

468   have been a better solution to develop a standardized dataset.

469   As genomic resources expand, conservation practitioners are faced with the problem of translating large

470   amounts of available data into useful management applications.  Challenges to enabling these applications

471   include reducing the number of markers to a level easily assayed in a cost-effective manner, working with

472   differing baseline datasets and missing loci, developing efficient analysis pipelines, and identifying

473   appropriate accuracy and precision for particular applications.  The method we employed in this study,

474   RADseq SNP discovery followed by GT-seq panel development, is a straightforward pathway to develop

475   genetic tools for fisheries management.  The panels we developed increased resolution of reporting groups in

476   Western Alaska Chinook salmon, a region with low population structure that has been historically difficult

477   for stock discrimination.  These panels offer a cost-effective method for improving genetic stock

478   identification for fisheries management. During the study, we also realized the difficulty in not only handling

479  large amounts of data with increasing numbers of SNPs but also leveraging data from varying laboratory

480  techniques.   We anticipate that  computational pipelines such as *GTscore* will expedite processing of both

481  RADseq, RAPTURE (Ali et al. 2016), and similar data sets and assist researchers in sharing and

482  standardizing data not only across techniques, but among researchers and laboratories.

483  **Acknowledgements**

497

**References**

Ali, O.A., O'Rourke, S.M., Amish, S.J., Meek, M.H., Luikart, G., Jeffres, C., and Miller, M.R. 2016. RAD capture (Rapture): flexible and efficient sequence-based genotyping. Genetics **202**(2): 389-400.

Allendorf, F.W., Hohenlohe, P.A., and Luikart, G. 2010. Genomics and the future of conservation genetics. Nature Reviews Genetics **11**(10): 697-709.

Anderson, E.C. 2010. Assessing the power of informative subsets of loci for population assignment: standard methods are upwardly biased. Molecular ecology resources **10**(4): 701-710.

Anderson, E.C., Waples, R.S., and Kalinowski, S.T. 2008. An improved method for predicting the accuracy of genetic stock identification. Canadian Journal of Fisheries and Aquatic Sciences **65**(7): 1475-1486.

Baetscher, D.S., Clemento, A.J., Ng, T.C., Anderson, E.C., and Garza, J.C. 2018. Microhaplotypes provide increased power from short-read DNA sequences for relationship inference. Molecular ecology resources **18**(2): 296-305.

Beacham, T.D., Wallace, C., MacConnachie, C., Jonsen, K., McIntosh, B., Candy, J.R., and Withler, R.E. 2018. Population and individual identification of Chinook salmon in British Columbia through parentage-based tagging and genetic stock identification with single nucleotide polymorphisms. Canadian Journal of Fisheries and Aquatic Sciences **75**(7): 1096-1105.

Beacham, T.D., Winter, I., Jonsen, K.L., Wetklo, M., Deng, L.T., and Candy, J.R. 2008. The application of rapid microsatellite-based stock identification to management of a Chinook salmon troll fishery off the Queen Charlotte Islands, British Columbia. North American Journal of Fisheries Management **28**(3): 849-855.

Bernatchez, L., Wellenreuther, M., Araneda, C., Ashton, D.T., Barth, J.M.I., Beacham, T.D., Maes, G.E., Martinsohn, J.T., Miller, K.M., Naish, K.A., Ovenden, J.R., Primmer, C.R., Suk, H.Y., Therkildsen, N.O., and Withler, R.E. 2017. Harnessing the Power of Genomics to Secure the Future of Seafood. Trends in Ecology & Evolution **32**(9): 665-680.

Bolstad, G.H., Hindar, K., Robertsen, G., Jonsson, B., Saegrov, H., Diserud, O.H., Fiske, P., Jensen, A.J., Urdal, K., Naesje, T.F., Barlaup, B.T., Floro-Larsen, B., Lo, H., Niemela, E., and Karlsson, S. 2017. Gene flow from domesticated escapes alters the life history of wild Atlantic salmon. Nature Ecology & Evolution **1**(5).

Campbell, N.R., Harmon, S.A., and Narum, S.R. 2015. Genotyping-in-Thousands by sequencing (GT-seq): a cost effective SNP genotyping method based on custom amplicon sequencing. Molecular ecology resources **15**(4): 855-867.

Candy, J.R., Campbell, N.R., Grinnell, M.H., Beacham, T.D., Larson, W.A., and Narum, S.R. 2015. Population differentiation determined from putative neutral and divergent adaptive genetic markers in Eulachon (*Thaleichthys pacificus*, *Osmeridae*), an anadromous Pacific smelt. Molecular ecology resources **15**(6): 1421-1434.

Chowen, T.R., and Nagler, J.J. 2005. Lack of sex specificity for growth hormone pseudogene in fall-run Chinook salmon from the Columbia River. Transactions of the American Fisheries Society **134**(1): 279-282.

24

530    Clemento, A.J., Crandall, E.D., Garza, J.C., and Anderson, E.C. 2014. Evaluation of a single nucleotide polymorphism

531    baseline for genetic stock identification of Chinook Salmon (Oncorhynchus tshawytscha) in the California Current

532    large marine ecosystem. Fishery Bulletin **112**(2-3): 112-130.

533    Dann, T.H., Habicht, C., Baker, T.T., and Seeb, J.E. 2013. Exploiting genetic diversity to balance conservation and

534    harvest of migratory salmon. Canadian Journal of Fisheries and Aquatic Sciences **70**(5): 785-793.

535    Dann, T.H., Habicht, C., Templin, W.D., Seeb, L.W., McKinney, G.J., and Seeb, J.S. 2018. Identification of genetic

536    markers useful for mixed stock analysis of Chinook salmon in Cook Inlet, Alaska. Alaska Department of Fish and

537    Game, Division of Commercial Fisheries, Anchorage.

538    deReynier, Y.L. 1998. Evolving principles of international fisheries law and the North Pacific anadromous fish

539    commission. Ocean Development and International Law **29**(2): 147-178.

540    Foll, M., and Gaggiotti, O. 2008. A genome-scan method to identify selected loci appropriate for both dominant and

541    codominant markers: a Bayesian perspective. Genetics **180**(2): 977-993.

542    Forseth, T., Barlaup, B.T., Finstad, B., Fiske, P., Gjoaester, H., Falkegard, M., Hindar, A., Mo, T.A., Rikardsen, A.H.,

543    Thorstad, E.B., Vollestad, L.A., and Wennevik, V. 2017. The major threats to Atlantic salmon in Norway. Ices Journal

544    of Marine Science **74**(6): 1496-1513.

545    Gilbey, J., Cauwelier, E., Coulson, M.W., Stradmeyer, L., Sampayo, J.N., Armstrong, A., Verspoor, E., Corrigan, L.,

546    Shelley, J., and Middlemas, S. 2016. Accuracy of Assignment of Atlantic Salmon (Salmo salar L.) to Rivers and

547    Regions in Scotland and Northeast England Based on Single Nucleotide Polymorphism (SNP) Markers. PloS one

548    **11**(10): e0164327.

549    Gilbey, J., Wennevik, V., Bradbury, I.R., Fiske, P., Hansen, L.P., Jacobsen, J.A., and Potter, T. 2017. Genetic stock

550    identification of Atlantic salmon caught in the Faroese fishery. Fisheries Research **187**: 110-119.

551    Gisclair, B.R. 2009. Salmon bycatch management in the Bering Sea walleye pollock fishery:  threats and

552    opportunityies for Western Alaska. *In* Pacific Salmon:  Ecology and Management of Western Alaska's Populations.

553    *Edited by* C.C. Krueger and C.E. Zimmerman. American Fisheries Society Symposium 70, Bethesda, Maryland. pp.

554    799-816.

555    Griffiths, J.R., Schindler, D.E., and Seeb, L.W. 2013. How Stock of Origin Affects Performance of Individuals across

556    a Meta-Ecosystem: An Example from Sockeye Salmon. PloS one **8**(3): e58584.

557    Healey, M.C. 1991. Life history of Chinook salmon. *In* Pacific Salmon Life Histories. *Edited by* C. Groot and L.

558    Margolis. UBC Press, Vancouver. pp. 311-394.

559    Hess, J.E., Campbell, N.R., Matala, A.P., Hasselman, D.J., and Narum, S.P. 2016. Genetic assessment of Columbia

560    River stocks, 4/1/2014-3/31/2105 annual report, 2008-907-00. CRITFC. Available from http://www.critfc.org/wp-

561    content/uploads/2016/04/16-03.pdf.

562  Holman, L.E., de la Serrana, D.G., Onoufriou, A., Hillestad, B., and Johnston, I.A. 2017. A workflow used to design

563  low density SNP panels for parentage assignment and traceability in aquaculture species and its validation in Atlantic

564  salmon. Aquaculture **476**: 59-64.

565  Langmead, B., and Salzberg, S.L. 2012. Fast gapped-read alignment with Bowtie 2. Nat Methods **9**(4): 357-359.

566  Larson, W.A., Seeb, J.E., Pascal, C.E., Templin, W.D., and Seeb, L.W. 2014a. Single-nucleotide polymorphisms

567  (SNPs) identified through genotyping-by-sequencing improve genetic stock identification of Chinook salmon

568  (*Oncorhynchus tshawytscha*) from western Alaska. Canadian Journal of Fisheries and Aquatic Sciences **71**(5): 698-

569  708.

570  Larson, W.A., Seeb, L.W., Everett, M.V., Waples, R.K., Templin, W.D., and Seeb, J.E. 2014b. Genotyping by

571  sequencing resolves shallow population structure to inform conservation of Chinook salmon (*Oncorhynchus

572  tshawytscha*). Evol Appl **7**(3): 355-369.

573  Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., and

574  Genome Project Data Processing, S. 2009. The Sequence Alignment/Map format and SAMtools. Bioinformatics

575  **25**(16): 2078-2079.

576  Limborg, M.T., Larson, W.A., Seeb, L.W., and Seeb, J.E. 2017. Screening of duplicated loci reveals hidden divergence

577  patterns in a complex salmonid genome. Molecular ecology.

578  Lin, P.C., Adams, R.M., and Berrens, R.P. 1996. Welfare effects of fishery policies: Native American treaty rights and

579  recreational salmon fishing. Journal of Agricultural and Resource Economics **21**(2): 263-276.

580  McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K., Altshuler, D.,

581  Gabriel, S., Daly, M., and DePristo, M.A. 2010. The Genome Analysis Toolkit: a MapReduce framework for

582  analyzing next-generation DNA sequencing data. Genome Res **20**(9): 1297-1303.

583  McKinney, G.J., Seeb, J.E., and Seeb, L.W. 2017. Managing mixed-stock fisheries: genotyping multi-SNP haplotypes

584  increases power for genetic stock identification. Canadian Journal of Fisheries and Aquatic Sciences **74**(4): 429-434.

585  McKinney, G.J., Seeb, L.W., Larson, W.A., Gomez-Uchida, D., Limborg, M.T., Brieuc, M.S., Everett, M.V., Naish,

586  K.A., Waples, R.K., and Seeb, J.E. 2016. An integrated linkage map reveals candidate genes underlying adaptive

587  variation in Chinook salmon (*Oncorhynchus tshawytscha*). Molecular ecology resources **16**(3): 769-783.

588  McKinney, G.J., Waples, R.K., Pascal, C.E., Seeb, L.W., and Seeb, J.E. 2018. Resolving allele dosage in duplicated

589  loci using genotyping-by-sequencing data: A path forward for population genetic analysis. Molecular ecology

590  resources **18**(3): 570-579.

591  Meek, M.H., Baerwald, M.R., Stephens, M.R., Goodbla, A., Miller, M.R., Tomalty, K.M.H., and May, B. 2016.

592  Sequencing improves our ability to study threatened migratory species: Genetic population assignment in California's

593  Central Valley Chinook salmon. Ecology and Evolution **6**(21): 7706-7716.

594  Miller, B.G. 1993. The press, the Boldt decision, and indian-white relations. American Indian Culture and Research

595  Journal **17**(2): 75-97.

596   Moran, B.M., and Anderson, E.C. 2018. Bayesian inference from the conditional genetic stock identification model.

597   Canadian Journal of Fisheries and Aquatic Sciences.

598   Myers, K.W., and Rogers, D.E. 1988. Stock Origins of Chinook Salmon in Incidental Catches by Groundfish Fisheries

599   in the Eastern Bering Sea. North American Journal of Fisheries Management **8**(2): 162-171.

600   Myers, K.W., Walker, R.V., Davis, N.D., Armstrong, J.L., and Kaeriyama, M. 2009. High seas distribution, biology,

601   and ecology of Arctic-Yukon-Kuskokwim salmon: direct information from high seas tagging experiments 1954-2006.

602   *In* Pacific Salmon: Ecology and Management of Western Alaska's Populations. *Edited by* C.C. Krueger and C.E.

603   Zimmerman. American Fisheries Society, Bethesda. pp. 201-240.

604   Nagler, J.J., Bouma, J., Thorgaard, G.H., and Dauble, D.D. 2001. High incidence of a male-specific genetic marker in

605   phenotypic female chinook salmon from the Columbia River. Environmental Health Perspectives **109**(1): 67-69.

606   Nagler, J.J., Cavileer, T., Steinhorst, K., and Devlin, R.H. 2004. Determination of genetic sex in chinook salmon

607   (Oncorhynchus tshawytscha) using the male-linked growth hormone pseudogene by real-time PCR. Marine

608   biotechnology **6**(2): 186-191.

609   Ohlberger, J., Ward, E.J., Schindler, D.E., and Lewis, B. 2018. Demographic changes in Chinook salmon across the

610   Northeast Pacific Ocean. Fish and Fisheries **19**(3): 533-546.

611   Phillips, R.B., Park, L.K., and Naish, K.A. 2013. Assignment of Chinook salmon (*Oncorhynchus tshawytscha*) linkage

612   groups to specific chromosomes reveals a karyotype with multiple rearrangements of the chromosome arms of rainbow

613   trout (*Oncorhynchus mykiss*). G3 **3**(12): 2289-2295.

614   Pritchard, V.L., Erkinaro, J., Kent, M.P., Niemelä, E., Orell, P., Lien, S., and Primmer, C.R. 2016. Single nucleotide

615   polymorphisms to discriminate different classes of hybrid between wild Atlantic salmon and aquaculture escapees.

616   Evolutionary Applications **9**(8): 1017-1031.

617   Rousset, F. 2008. Genepop'007: A complete re-implementation of the Genepop software for Windows and Linux.

618   Molecular ecology resources **8**(1): 103-106.

619   Schoen, E.R., Wipfli, M.S., Trammell, E.J., Rinella, D.J., Floyd, A.L., Grunblatt, J., McCarthy, M.D., Meyer, B.E.,

620   Morton, J.M., Powell, J.E., Prakash, A., Reimer, M.N., Stuefer, S.L., Toniolo, H., Wells, B.M., and Witmer, F.D.W.

621   2017. Future of Pacific Salmon in the Face of Environmental Change: Lessons from One of the World's Remaining

622   Productive Salmon Regions. Fisheries **42**(10): 538-+.

623   Siegel, J.E., McPhee, M.V., and Adkison, M.D. 2017. Evidence that Marine Temperatures Influence Growth and

624   Maturation of Western Alaskan Chinook Salmon. Marine and Coastal Fisheries **9**(1): 441-456.

625   Smit, A., Hubley, R., and Green, P. 2013. RepeatMasker Open-4.0. available at: http://repeatmasker.org.

626   Smith, C.T., Elfstrom, C.M., Seeb, L.W., and Seeb, J.E. 2005a. Use of sequence data from rainbow trout and Atlantic

627   salmon for SNP detection in Pacific salmon. Molecular Ecology **14**(13): 4193-4203.

628   Smith, C.T., Seeb, J.E., Schwenke, P., and Seeb, L.W. 2005b. Use of the 5 '-nuclease reaction for single nucleotide

629   polymorphism genotyping in Chinook salmon. Transactions of the American Fisheries Society **134**(1): 207-217.

27

630     Smith, C.T., Templin, W.D., Seeb, J.E., and Seeb, L.W. 2005c. Single Nucleotide Polymorphisms (SNPs) provide

631     rapid and accurate estimates of the proportions of U.S. and Canadian Chinook salmon caught in Yukon River fisheries.

632     North American Journal of Fisheries Management **25**(3): 944-953.

633     Sylvester, E.V.A., Bentzen, P., Bradbury, I.R., Clément, M., Pearce, J., Horne, J., and Beiko, R.G. 2017. Applications

634     of random forest feature selection for fine-scale genetic population assignment. Evolutionary Applications: n/a-n/a.

635     Templin, W.D., Seeb, J.E., Jasper, J.R., Barclay, A.W., and Seeb, L.W. 2011. Genetic differentiation of Alaska

636     Chinook salmon: the missing link for migratory studies. Molecular ecology resources **11 Suppl 1**: 226-246.

637     Thorgaard, G.H. 1977. HETEROMORPHIC SEX-CHROMOSOMES IN MALE RAINBOW-TROUT. Science

638     **196**(4292): 900-902.

639     Von Bargen, J., Smith, C.T., and Reuth, J. 2015. Development of a Chinook salmon sex identification SNP assay

640     based on the growth hormone pseudogene. Journal of Fish and Wildlife Management **6**(1): 213-219.

641     Walsh, V.M. 1998. Eliminating driftnets from the North Pacific Ocean: US-Japanese cooperation in the International

642     North Pacific Fisheries Commission, 1953-1993. Ocean Development and International Law **29**(4): 295-322.

643     Waples, R.K., Seeb, J.E., and Seeb, L.W. 2017. Congruent population structure across paralogous and nonparalogous

644     loci in Salish Sea chum salmon (*Oncorhynchus keta*). Molecular ecology **26**(16): 4131-4144.

645     Warheit, K.I., Seeb, L., Templin, W.D., and Seeb, J. 2013. Moving GSI into the next decade:  SNP coordination for

646     Pacific Salmon Treaty fisheries.  Washington Department of Fish and Wildlife,  Report FPT 13-09.

647     http://wdfw.wa.gov/publications/01629/wdfw01629.pdf.

648     You, F.M., Huo, N., Gu, Y.Q., Luo, M.C., Ma, Y., Hane, D., Lazo, G.R., Dvorak, J., and Anderson, O.D. 2008.

649     BatchPrimer3: a high throughput web application for PCR and sequencing primer design. BMC Bioinformatics **9**: 253.

651    **Table 1**.  Number of samples sequenced and retained per population for RADseq and GT-seq analyses.

| Population | Region | RADseq | | GT-seq | | Total | Source[1] |
| | | Sequenced | Retained | Sequenced | Retained | | |
|---|---|---|---|---|---|---|---|
| Pitka Fork | Upper Kuskokwim River | 0 | 0 | 95 | 95 | 95 | 1 |
| Takotna | Upper Kuskokwim River | 0 | 0 | 95 | 94 | 94 | 1 |
| Tatlawiksuk | Upper Kuskokwim River | 0 | 0 | 95 | 56 | 56 | 1 |
| Necons | Upper Kuskokwim River | 48 | 47 | 48 | 48 | 95 | 2 |
| George | Kuskokwim River | 48 | 32 | 49 | 40 | 72 | 2 |
| Kogrukluk | Kuskokwim River | 64 | 48 | 47 | 47 | 95 | 3 |
| Aniak | Kuskokwim River | 48 | 47 | 48 | 48 | 95 | 1 |
| Kisaralik | Kuskokwim River | 48 | 48 | 47 | 47 | 95 | 1 |
| Kwethluk | Kuskokwim River | 48 | 35 | 52 | 52 | 90 | 1 |
| Eek | Kuskokwim River | 0 | 0 | 95 | 93 | 93 | 1 |
| Kanektok | Kuskokwim Bay | 48 | 44 | 49 | 49 | 95 | 1 |
| Arolik | Kuskokwim Bay | 48 | 46 | 47 | 47 | 93 | 1 |
| Goodnews | Goodnews Bay | 48 | 47 | 48 | 48 | 95 | 2 |
| Togiak | Togiak Bay | 48 | 46 | 47 | 45 | 91 | 1 |
| Iowithla | Nushagak River | 48 | 47 | 47 | 17 | 64 | 1 |
| Stuyahok | Nushagak River | 48 | 48 | 47 | 45 | 93 | 1 |
| Koktuli | Nushagak River | 56 | 56 | 39 | 39 | 95 | 3 |
| Total | | 648 | 591 | 995 | 910 | 1506 | |

652    [1] Sources:  1) this study, 2) McKinney et al. (2018), 3)  Larson et al. (2014b).

**Table 2**. Putative sex loci identified by the presence of a two genotype classes (heterozygote and one homozygote) at approximately even frequencies in the population. Location of each locus on the Chinook salmon linkage map is listed for mapped loci. Accuracy was calculated by comparing genotypes with observed sex of 23 individuals from the Togiak River population.

| Locus | Chromosome | Position (cM) | Accuracy |
|---|---|---|---|
| *RAD67724* | *Ots17* | 4.7 | 83% |
| *RAD29719* | *Ots17* | 9.4 | 78% |
| *RAD27492* | NA | NA | 91% |

30

**Table 3**. Size of Chinook salmon linkage groups and number of loci for Washington linkage map, Alaska linkage map, and combined linkage map.

| Linkage Group | Washington Map Size (cM) | Singleton Loci | Duplicated Loci | Total Loci | Alaska Map Size (cM) | Singleton Loci | Duplicated Loci | Total Loci | Combined Map Size (cM) | Singleton Loci | Duplicated Loci | Total Loci |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ots01 | 125.85 | 669 | 103 | 772 | 107.72 | 674 | 107 | 781 | 117.24 | 1066 | 159 | 1225 |
| Ots02 | 120.4 | 428 | 140 | 568 | 119.94 | 491 | 155 | 646 | 118.78 | 729 | 231 | 960 |
| Ots03 | 129.64 | 507 | 140 | 647 | 109.09 | 517 | 161 | 678 | 125.78 | 788 | 236 | 1024 |
| Ots04 | 126.21 | 417 | 168 | 585 | 104.63 | 434 | 196 | 630 | 112.29 | 681 | 275 | 956 |
| Ots05 | 116.49 | 499 | 15 | 514 | 107.42 | 592 | 19 | 611 | 107.80 | 868 | 31 | 899 |
| Ots06 | 119.09 | 604 | 116 | 720 | 98.14 | 650 | 118 | 768 | 108.89 | 958 | 175 | 1133 |
| Ots07 | 134.24 | 538 | 113 | 651 | 110.29 | 577 | 157 | 734 | 120.00 | 850 | 216 | 1066 |
| Ots08 | 106.38 | 584 | 11 | 595 | 103.54 | 614 | 18 | 632 | 109.02 | 921 | 30 | 951 |
| Ots09 | 126.14 | 560 | 183 | 743 | 127.46 | 590 | 205 | 795 | 116.79 | 890 | 295 | 1185 |
| Ots10 | 122.51 | 410 | 23 | 433 | 106.76 | 504 | 22 | 526 | 111.70 | 739 | 37 | 776 |
| Ots11 | 103.84 | 378 | 53 | 431 | 108.11 | 432 | 137 | 569 | 104.64 | 627 | 168 | 795 |
| Ots12 | 127.67 | 462 | 159 | 621 | 102.06 | 511 | 195 | 706 | 120.69 | 780 | 268 | 1048 |
| Ots13 | 114.07 | 574 | 21 | 595 | 110.24 | 596 | 29 | 625 | 111.91 | 901 | 43 | 944 |
| Ots14 | 123.57 | 374 | 123 | 497 | 105.84 | 396 | 154 | 550 | 118.43 | 610 | 208 | 818 |
| Ots15 | 96.29 | 295 | 139 | 434 | 107.36 | 288 | 156 | 444 | 98.24 | 453 | 232 | 685 |
| Ots16 | 114.35 | 414 | 13 | 427 | 116.55 | 469 | 15 | 484 | 115.76 | 702 | 23 | 725 |
| Ots17 | 66.53 | 164 | 129 | 293 | 64.20 | 166 | 122 | 288 | 70.34 | 260 | 194 | 454 |
| Ots18 | 66.47 | 277 | 15 | 292 | 53.51 | 305 | 13 | 318 | 60.63 | 464 | 20 | 484 |
| Ots19 | 73.98 | 445 | 18 | 463 | 85.89 | 449 | 14 | 463 | 79.15 | 677 | 25 | 702 |
| Ots20 | 74.7 | 354 | 10 | 364 | 77.28 | 374 | 19 | 393 | 79.07 | 556 | 23 | 579 |
| Ots21 | 70.02 | 223 | 12 | 235 | 52.24 | 260 | 8 | 268 | 60.58 | 391 | 18 | 409 |
| Ots22 | 66.77 | 306 | 11 | 317 | 72.21 | 292 | 14 | 306 | 66.90 | 444 | 25 | 469 |
| Ots23 | 60.51 | 131 | 148 | 279 | 64.51 | 156 | 148 | 304 | 60.35 | 235 | 218 | 453 |

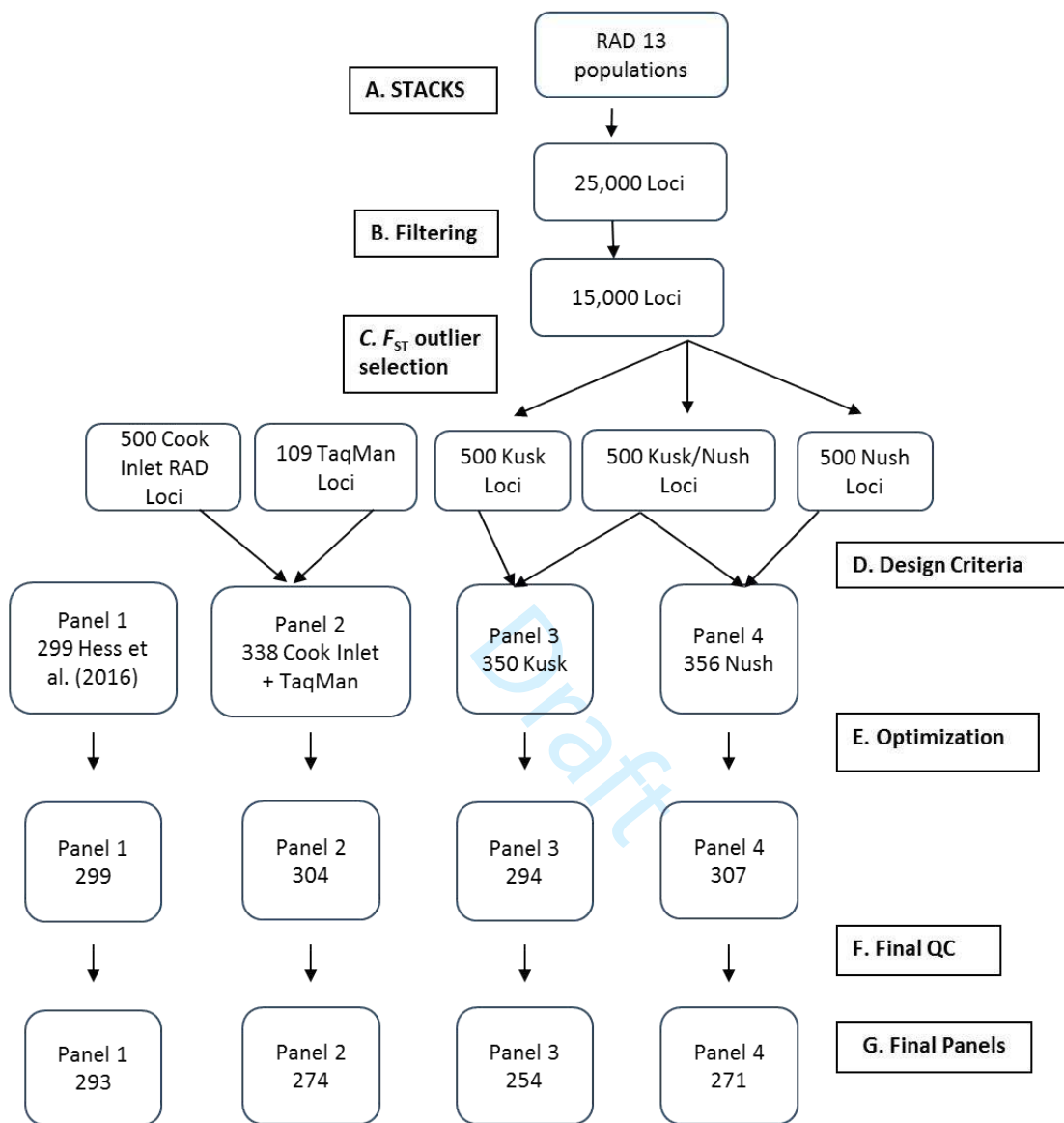| Linkage Group | Washington Map | | | | Alaska Map | | | | Combined Map | | | |
| | Size (cM) | Singleton Loci | Duplicated Loci | Total Loci | Size (cM) | Singleton Loci | Duplicated Loci | Total Loci | Size (cM) | Singleton Loci | Duplicated Loci | Total Loci |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ots24 | 56.61 | 217 | 2 | 219 | 55.31 | 245 | 9 | 254 | 56.30 | 333 | 11 | 344 |
| Ots25 | 73.55 | 303 | 14 | 317 | 70.35 | 311 | 16 | 327 | 69.16 | 482 | 25 | 507 |
| Ots26 | 66.66 | 342 | 13 | 355 | 56.62 | 332 | 15 | 347 | 66.16 | 522 | 24 | 546 |
| Ots27 | 69.63 | 145 | 150 | 295 | 56.01 | 152 | 166 | 318 | 62.69 | 231 | 250 | 481 |
| Ots28 | 64.39 | 280 | 10 | 290 | 55.12 | 323 | 8 | 331 | 65.15 | 467 | 14 | 481 |
| Ots29 | 57.75 | 234 | 6 | 240 | 51.66 | 265 | 8 | 273 | 55.40 | 378 | 11 | 389 |
| Ots30 | 70.21 | 343 | 12 | 355 | 58.71 | 342 | 7 | 349 | 63.50 | 536 | 18 | 554 |
| Ots31 | 48.8 | 239 | 9 | 248 | 55.18 | 227 | 6 | 233 | 50.89 | 365 | 13 | 378 |
| Ots32 | 68.42 | 122 | 139 | 261 | 52.10 | 110 | 159 | 269 | 61.98 | 182 | 237 | 419 |
| Ots33 | 77.29 | 335 | 3 | 338 | 88.48 | 337 | 6 | 343 | 79.96 | 508 | 9 | 517 |
| Ots34 | 80.74 | 111 | 115 | 226 | 59.49 | 103 | 132 | 235 | 77.19 | 168 | 191 | 359 |
| | | | | | | | | | | | | |
| Total | 3119.77 | 12,284 | 2,336 | 14,620 | 2874.02 | 13,084 | 2,714 | 15,798 | 3003.36 | 19,762 | 3,953 | 23,715 |

32

**Table 4.** Results of locus filtering prior to mixed-stock analysis (MSA). The number of total loci for each panel are in the Panel Loci column; the number of loci retained after filtering are in the *GSIsim* Loci column. The number of loci filtered at each step are listed in the remaining columns. The retained 847 singleton loci were evaluated for MSA using *GSIsim*.

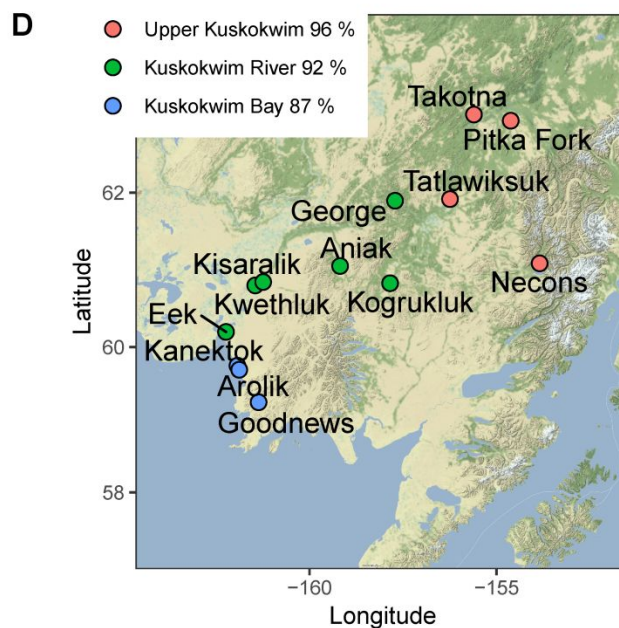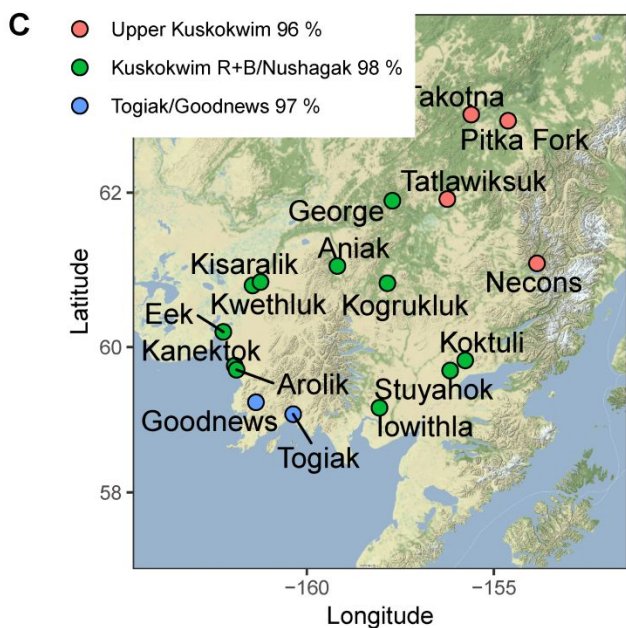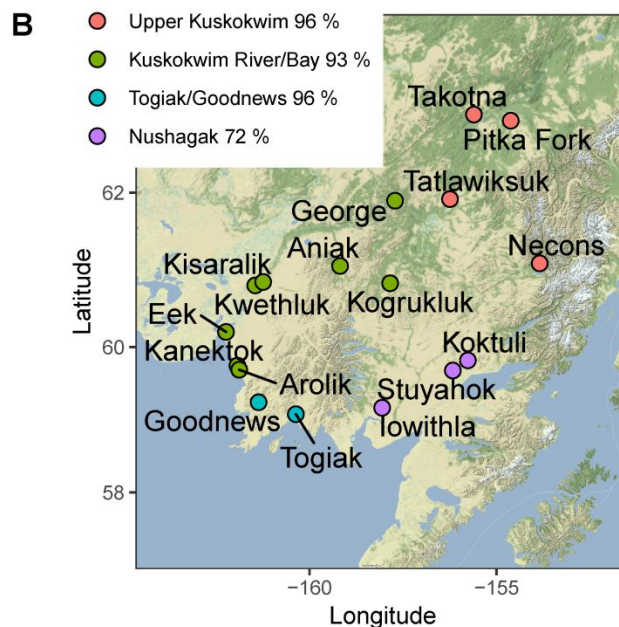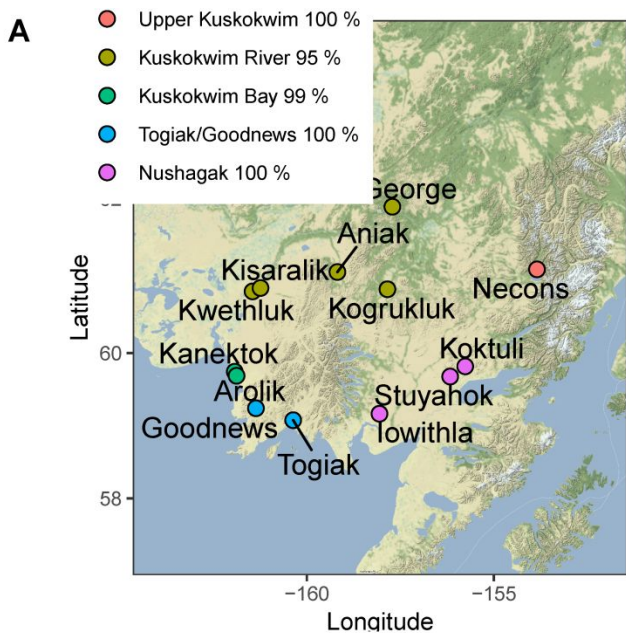| Panel | Panel Number | Panel Loci | No data for RADseq samples | Duplicate | Allele Frequency Discrepancies | Genotype Rate <70% | Mono-morphic | *GSIsim* Loci |
|---|---|---|---|---|---|---|---|---|
| Hess et al. (2016) | 1 | 293 | 61 | 6 | 2 | 33 | 34 | 157 |
| Cook Inlet + TaqMan | 2 | 274 | 23 | 13 | 4 | 1 | 4 | 232 |
| Kuskokwim | 3 | 254 | 0 | 24 | 6 | 0 | 2 | 219 |
| Nushagak | 4 | 271 | 0 | 27 | 5 | 0 | 0 | 239 |
| Total | | 1,092 | 84 | 70 | 17 | 34 | 40 | 847 |

**Figure 1.** Sample sites for populations used in this study. Shapes denote type of sequencing done on the population; populations are color coded by region. Map tiles by Stamen Design, under CC BY 3.0.

.

A. STACKS

RAD 13 populations

25,000 Loci

B. Filtering

15,000 Loci

C. $F_{ST}$ outlier selection

500 Cook Inlet RAD Loci

109 TaqMan Loci

500 Kusk Loci

500 Kusk/Nush Loci

500 Nush Loci

D. Design Criteria

Panel 1
299 Hess et al. (2016)

Panel 2
338 Cook Inlet + TaqMan

Panel 3
350 Kusk

Panel 4
356 Nush

E. Optimization

Panel 1
299

Panel 2
304

Panel 3
294

Panel 4
307

F. Final QC

Panel 1
293

Panel 2
274
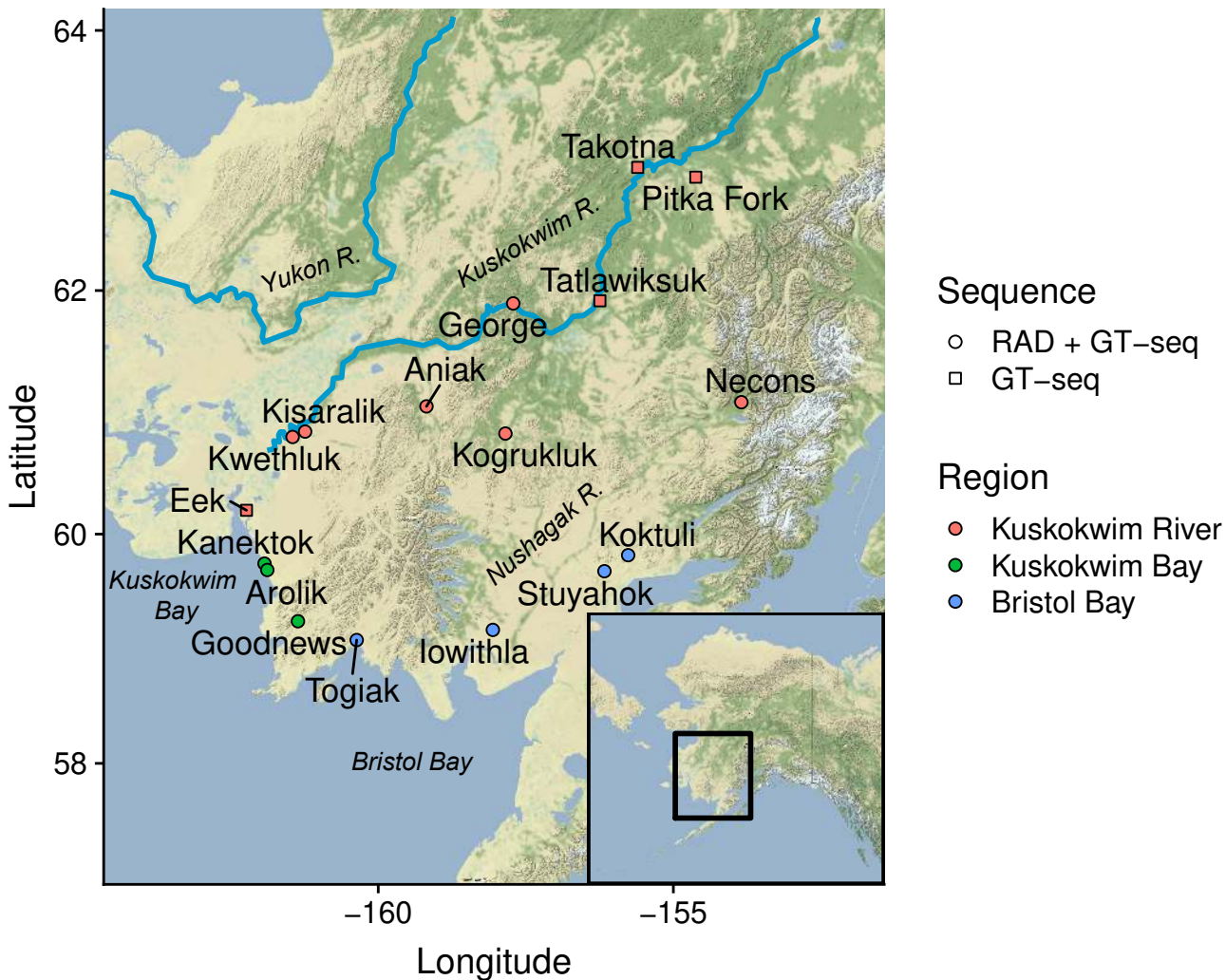
Panel 3
254

Panel 4
271

G. Final Panels

**Figure 2**.  Flowchart of GT-seq panel development.

**Figure 3.** Map of populations used in this study and reporting groups tested. SNP discovery was conducted on ascertainment populations using RADseq. All RADseq populations were used for genetic stock assignment along with additional populations and individuals genotyped using GT-seq. For each reporting group scenario, populations on the map are color coded by reporting group, and mean accuracy of reporting groups are shown in the legend. The maximum possible resolution using the full RADseq dataset (>15,000 loci) is shown in A. Reporting group scenarios for GT-seq panel (847 loci) resolution are B) fine-scale resolution, C) broad-scale resolution, D) resolution within Kuskokwim Bay/River, and E) resolution within Togiak Bay/Nushagak River. Map tiles by Stamen Design, under CC BY 3.0.
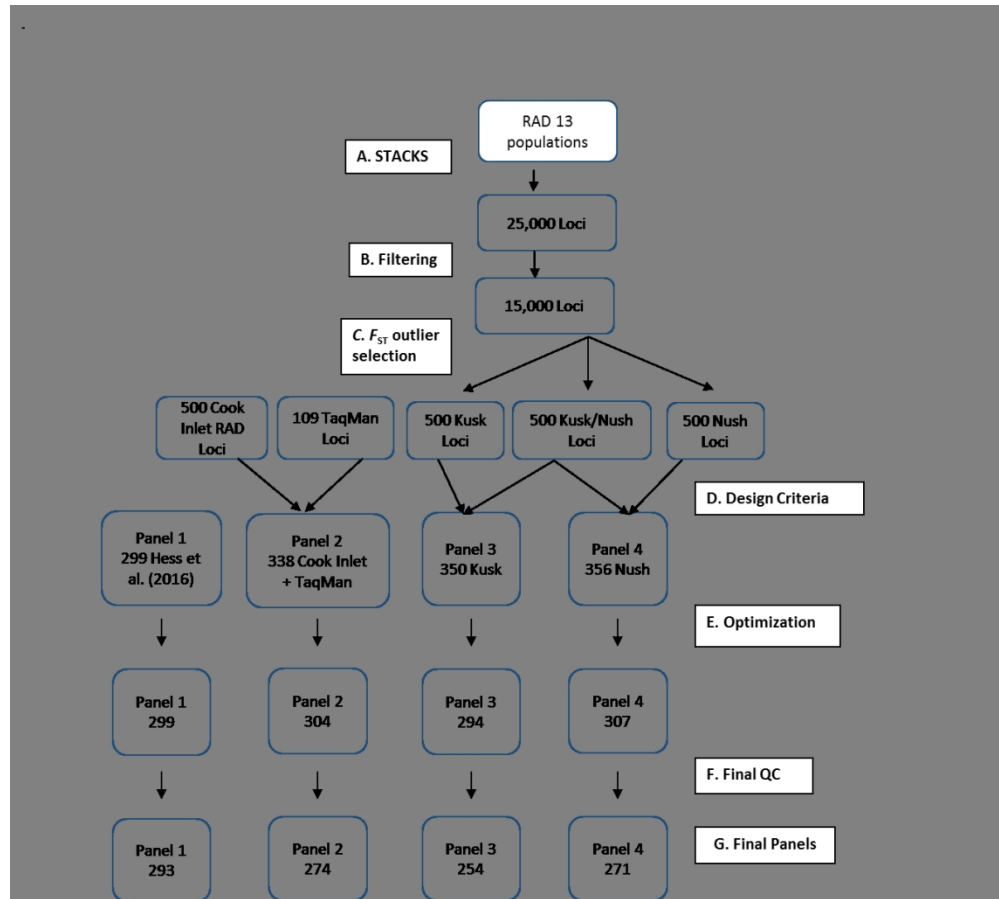
Figure 2. Flowchart of GT-seq panel development.

**A**
- Upper Kuskokwim 100 %
- Kuskokwim River 95 %
- Kuskokwim Bay 99 %
- Togiak/Goodnews 100 %
- Nushagak 100 %

**B**
- Upper Kuskokwim 96 %
- Kuskokwim River/Bay 93 %
- Togiak/Goodnews 96 %
- Nushagak 72 %

**C**
- Upper Kuskokwim 96 %
- Kuskokwim R+B/Nushagak 98 %
- Togiak/Goodnews 97 %

**D**
- Upper Kuskokwim 96 %
- Kuskokwim River 92 %
- Kuskokwim Bay 87 %

**E**
- Togiak 100 %
- Lower Nushagak 99 %
- Middle Nushagak 99 %