1 # Isolation-by-distance and population-size history inferences from

2 # the coho salmon (*Oncorhynchus kisutch*) genome

3

4 Eric B. Rondeau[1,2,3*], Kris A. Christensen[1,2*], David R. Minkley[1,2], Jong S. Leong[1], Michelle

5 T.T. Chan[2,4], Cody A. Despins[1], Anita Mueller[1], Dionne Sakhrani[2], Carlo A. Biagi[2], Quentin

6 Rougemont[5,6], Eric Normandeau[5], Steven J.M. Jones[7], Robert H. Devlin[2], Ruth E.

7 Withler[3], Terry D. Beacham[3], Kerry A. Naish[8], José M. Yáñez[9,11], Roberto Neira[10,11], Louis

8 Bernatchez[5], William S. Davidson[4], Ben F. Koop[1].

9

10 **Affiliations**

11 1: Department of Biology, Centre for Biomedical Research, University of Victoria,

12 Victoria, BC, V8W 2Y2, Canada

13 2: Fisheries and Oceans Canada, 4160 Marine Drive, West Vancouver, BC, V7V 1N6,

14 Canada

15 3: Fisheries and Oceans Canada, Pacific Biological Station, 3190 Hammond Bay Road,

16 Nanaimo, BC, V9T 6N7, Canada

17 4: Department of Molecular Biology and Biochemistry, Simon Fraser University,

18 Burnaby, V5A 1S6, Canada

19 5: Institut de Biologie Intégrative et des Systèmes (IBIS), Université Laval, Québec, QC,

20 G1V 0A6, Canada

21 6: Current: CEFE, Univ Montpellier, CNRS, EPHE, IRD, Montpellier, France

22 7: Canada's Michael Smith Genome Sciences Centre, BC Cancer, Vancouver, British

23 Columbia, Canada

24 8: School of Aquatic and Fishery Sciences, University of Washington, Seattle, WA 98105,

25 USA

26 9: Facultad de Ciencias Veterinarias y Pecuarias, Universidad de Chile, Santa Rosa 11735,

27 La Pintana, Santiago, Chile

28    10: Facultad de Ciencias Agronómicas, Universidad de Chile, Santa Rosa 11315, La

29    Pintana, Santiago, Chile

30    11: Millennium Nucleus of Austral Invasive Salmonids (INVASAL), Concepción, Chile

31    *equal contributions

32

33    **Corresponding authors**

34    Ben F. Koop: bkoop@uvic.ca

35    Kris A. Christensen: kris.christensen@wsu.edu

36

37    **Running Head**

38    coho salmon genome assembly

39

43

44    **Abstract**

45        Coho salmon (*Oncorhynchus kisutch*) are a culturally and economically important

46    species that return from multiyear ocean migrations to spawn in rivers that flow to the

47    Northern Pacific Ocean. Southern stocks of coho salmon have significantly declined over

48    the past quarter century, and unfortunately, conservation efforts have not reversed this

49    trend. To assist in stock management and conservation efforts, we generated two

50    chromosome-level genome assemblies and sequenced 24 RNA-seq libraries to better

51    annotate the coho salmon genome assemblies. We also resequenced the genomes of 83

52    coho salmon across their North American range to identify nucleotide variants,

53    characterize the broad effects of isolation-by-distance using a genome-wide association

54    analysis approach, and understand the demographic histories of these salmon by

55    modeling population size from genome-wide data. We observed that more than 13% of

2

56  all SNPs were associated with latitude (before multiple test correction), likely an affect

57  of isolation-by-distance. From demographic history modeling, we estimated that the

58  SNP latitudinal gradient likely developed as recently as 8,000 years ago. In addition, we

59  identified four genes each harboring multiple SNPs associated with latitude; all of these

60  SNPs were also predicted to modify the function of the gene. Three of these genes have

61  roles in cell junction maintenance and may be involved in osmoregulation. This signifies

62  that ocean salinity may have been a factor influencing coho salmon recolonization after

63  the last glaciation period – generating the current pattern of variation in these three

64  genes.

65

66  **Introduction**

67      Coho salmon have special cultural significance to the people of the First Nations

68  in British Columbia and have traditionally been one of the highest-value Pacific salmon

69  in the commercial and recreational fishery sectors. In 1977, a climatic regime shift in the

70  North Pacific Ocean ushered in three decades of increasing global salmon production

71  that culminated in 2009, when over 600 million salmon (1.1 million metric tonnes) were

72  harvested [1]. However, this increased production of salmon masked substantial

73  variability in regional abundances and species composition. Whereas the productivity

74  and harvest of chum (*Oncorhynchus keta*), pink (*Oncorhynchus gorbuscha*), and sockeye

75  salmon (*Oncorhynchus nerka*) increased throughout the North Pacific after 1977, the

76  opposite was true for coho and Chinook salmon (*Oncorhynchus tshawytscha*). These

77  declines became particularly acute after 1989 when marine survival for these species

78  began a downward spiral that has yet to be reversed [1, 2]. A severe decline in the

79  highly lucrative recreational coho salmon fishery in the Strait of Georgia saw the

80  numbers of fish caught decline from an average of over 500,000 to less than 100,000

81  throughout the 1990s [3]. In 2004, the recreational catch in the Strait of Georgia was

82  9,500 coho salmon [4].

83      In British Columbia (BC), the Salmon Enhancement Program (SEP) was launched

3

84    to double salmon production with the establishment of 18 major Department of

85    Fisheries and Oceans (DFO) operated hatchery facilities and spawning channels.

86    Throughout a stable harvest period of the 1980s, SEP releases of coho salmon juveniles

87    increased from 5 to over 20 million, and the proportion of hatchery salmon in the

88    fisheries increased from 5 to 20%. A precipitous decline in coho salmon production

89    occurred in the 1990s. The decline was even more dramatic in the inner waters of the

90    Strait of Georgia and Puget Sound. By the end of the decade, the commercial coho

91    salmon fishery was closed and 'marked or hatchery-only' recreational fisheries had been

92    instituted as a wild coho salmon conservation measure in southern BC.

93    The concern that hatchery fish were replacing wild fish was raised and, indeed,

94    by 1998 70% of coho salmon in the Strait of Georgia were of hatchery origin [5]. From

95    1998 to 2007, the survival of the Strait of Georgia coho salmon over the first four

96    months after entering a marine environment (May-September) decreased from 15% to

97    1% [2]. Processes associated with the low early marine survival remains unknown, but

98    marine climatic changes were implicated and hatchery salmon survival was even lower

99    than for wild salmon [2]. These results led to renewed calls for improved strategies for

100    wild and hatchery coho salmon management, and a re-evaluation of wild-hatchery

101    interactions in the species [2]. As such, a call was made to understand genetic influences

102    on coho salmon survival and to produce high-quality genomic resources such as a

103    chromosome-level reference genome assembly to enable technological support in

104    informing management practices and decision within this species.

105    In a large-scale coho salmon population structure analyses of coho salmon

106    sampled from 318 localities, in 38 different regional groups in North America and Russia

107    (representing most of the natural distribution of coho salmon), 17 microsatellite loci

108    showed that salmon clustered geographically and regions could be delineated along a

109    north – south gradient, with reduced variation to the north and isolated inland

110    populations [6]. These results were refined with increased genetic markers, finding that

111    isolation-by-distance from a main southern glacial refugia after the last ice-age could

4

112    explain most of the patterns of genetic diversity in modern coho salmon across the

113    North American distribution [7, 8]. These last two studies were supported by the

114    reference genome assemblies described in this study and illustrate how important such

115    resources are for understanding the basic biology of a species.

116         With this in mind, the goals of this study were to expand upon our basic

117    understanding of the coho salmon genome and help build upon the knowledge of the

118    already excellent framework of population structure mentioned above. Our method to

119    do this was to construct a high quality, annotated reference genome assembly and by

120    building a comprehensive inventory of genetic variation (SNPs) from a wide

121    geographical distribution. From the complete SNP dataset, we were then able to expand

122    upon what was known about isolation-by-distance and demographic history of coho

123    salmon. RNA-seq data for various tissues were also generated to facilitate genome

124    annotation by the NCBI.

125

126

127    **Materials and Methods**

128    *Coho salmon samples for genome assembly*

129         All animals were reared in compliance with the Canadian Council on Animal Care

130    Guidelines, under permit from the Fisheries and Oceans Canada Pacific Region Animal

131    Care Committee (under Ex.7.1). Using Inch Creek coho salmon, we generated fully

132    homozygous diploid gynogenetic individuals (doubled haploids) to help improve genome

133    assembly quality (as noted in [9]). For details on doubled haploid generation and DNA

134    extraction methods, please see the Supplemental Methods section. Tissues were also

135    collected for RNA-seq and included kidney, heart, head kidney, spleen, gill, nares, ovary,

136    white muscle, brain, eye, gut, liver, skin, stomach, and pyloric caecum. See the

137    Supplemental Methods for further details on RNA extraction.

138

139    *Genome sequence and assembly – Version 1*

140         A common sequencing and assembly pipeline for salmonids was used for this

141    version of the genome assembly (e.g., [10–12]). Full details of sequencing and genome

142    assembly can be found in the Supplemental Methods section. A brief description of the

143    assembly involved generating Illumina libraries (mate-pair and paired end), generating

144    PacBio data, and assembling the Illumina sequence data using Allpaths-LG [13] followed

145    by scaffolding with PacBio data using PBJelly [14]. All sequencing data related to this

146    genome assembly and annotation were submitted to the NCBI under BioProject

147    Accession: PRJNA352719. Genome completeness was assessed using BUSCO (v3.0) [15],

148    with default settings aside from "-sp zebrafish" using the ODB9 Actinopterygian

149    database.

150         A circos plot (v0.69-4) was generated to show the relationship of homeologous

151    chromosome resulting from the salmonid specific genome duplication [16, 17] for both

152    versions of the genome assembly. For further details of the circos plot see Supplemental

153    Methods.

154

155    *Genome sequence and assembly – Version 2*

156         Version 2 of the coho salmon genome assembly incorporated 10X chromium

157    data, Hi-C data, and new PacBio data with the previous Illumina sequencing and PacBio

158    data generated for the first version. Some of this data came from a different doubled-

159    haploid coho salmon individual compared to the first version. For full details of

160    sequencing and assembly methodology see Supplemental Methods.

161

162    *Transciptome and annotation*

163         RNA-seq data was generated from 15 tissues taken from the same doubled-

164    haploid coho salmon used to produce the first genome assembly. RNA-seq data was also

165    generated from two other coho salmon for this project, including: spleen, head kidney,

166    kidney, gill, and gut (gut was only from one of the two salmon). In total, RNA from 24

167    tissues were sent for library construction and sequencing at the Michael Smith Genome

6

168    Sciences Centre (Vancouver, BC, Canada). Eukaryotic single-strand RNAseq libraries

169    were prepared and sequenced across 7 lanes of PE125 sequencing on an Illumina

170    HiSeq2500. Four tissues were pooled per lane except brain, ovary, liver and gut from the

171    genome individual (DH3), which were pooled two per lane. Sequences were submitted

172    to NCBI under SRR5333359-SRR5333382 for eventual inclusion in the standard NCBI

173    Eukaryotic Genome Annotation pipeline, which has been used on many genome

174    assemblies. This data was generated for use in the NCBI annotation (Version 2), but was

175    not used in any other way in this study.

176

177    *Repeat library*

178        A species-specific repeat library was generated for coho salmon using the

179    methodology developed for salmonids in [18], and fully described in [10]. In brief, the

180    Atlantic salmon repeat library [18], was combined with repetitive sequences from the

181    RepBase database [19]. The RepBase sequences were derived from the Salmoniformes

182    family. They excluded simple repeats. RepeatModeler v1.0.8 [20] was also used

183    together with the genome assembly in a *de novo* approach. The repetitive sequences

184    were then aligned to the coho genome with BLASTN [21]. Sequences were classified into

185    either high-confidence or low-confidence categories based on frequency and length.

186    Low-confidence repeats were removed, and after filtering all of the sequences were

187    compared to each other using an all-by-all BLASTN search. A redundancy filter was

188    applied, prioritizing longest and highest-confidence repeats where two sequences were

189    considered to overlap.

190

191    *Whole-genome resequencing and nucleotide variant calling*

192        Whole-genome resequencing was used to characterize broad genomic

193    characteristics across the coho salmon's North American range. Table 1 contains a list of

194    sampled locations (see File S1 for more information). We included one commercial

195    strain from Chile as well (Table 1).

7

196       DNA was extracted from fin-clips using the DNeasy Blood and Tissue extraction

197      kit (Qiagen) or a MagMAX DNA Multi-Sample Ultra Kit with a KingFisher (ThermoFisher

198      Scientific). Following DNA extraction, samples were quantified by Qubit BR DNA assay

199      (ThermoFisher) and integrity validated by agarose gel electrophoresis. At McGill

200      University and Genome Québec Innovation Centre (Montreal, QC, Canada), individual

201      Illumina libraries were constructed with Illumina TruSeq LT sample preparation kits, and

202      each individual was sequenced separately on a lane of Illumina HiSeq2500 PE125 or in

203      batches of four on a HiSeqXTen (PE150) lane, targeting approximately 15-30X coverage.

204      Resequenced genomes were submitted to the NCBI under BioProject:PRJNA401427 and

205      PRJNA808051 (File S1).

206       Nucleotide variant calling on the dataset followed GATK3 best practices where

207      possible. BWA-MEM v0.7.17 [22] was used to align Illumina data to the reference

208      genome (version 2), with -M option for Picard compatibility. The Picard v2.18.9 [23]

209      AddOrReplaceReadGroups program was used to add read group IDs, and the

210      MarkDuplicates program was used to mark duplicates (default settings). GATK v3.8 [24,

211      25] was then used to call genotypes. Base and variant recalibration were each

212      performed once (for two rounds through genotyper). The variants used for recalibration

213      were from 1) a reduced set of very high-confidence calls following default "hard-

214      filtering" guidelines from GATK documentation from the first round of genotyping with a

215      particular focus on coding regions, and 2) validated SNPs on a 200K Affymetrix SNParray

216      [26].

217       Following genotyping, VCFtools v0.1.15 [27] was used to additionally thin data to

218      only include biallelic SNPs with a minor allele frequency of 0.05 or greater, variants with

219      fewer than 10% missing genotypes, and variants with a mean coverage between 5 and

220      200. Minor allele frequency was not used for filtering in the SMC++ analysis (below).

221      Some individuals were removed at this point from the VCF file because they were not

222      intended for this study (see NCBI BioProject: PRJNA808051 and File S1 for removed

223      samples). They were included as it was more computationally efficient to call all

8

224    individuals at the same time. Finally, the SNPs were filtered for linkage disequilibrium to

225    reduce the influence of large haploblocks in the principal components analysis (PCA)

226    (bcftools [28] version 1.9-102-g958180e +prune -w 20kb -l 0.4 -n 2).

227

228    *Whole-genome analyses*

229         A PCA was performed with variants that had been filtered for linkage

230    disequilibrium (see previous paragraph) using PLINK [29, 30] v1.90b6.15 with default

231    parameters. PLINK was also used to identify and quantify runs of homozygosity using

232    default settings (Figure S1). For comparison, we also performed the analysis with the

233    following parameters (Figure S2): min SNP count – 100, min length – 100 kb, max

234    inverse density – 50 kb/SNP, max internal gap 100 kb, max heterozygous genotypes 1,

235    SNP scanning window size – 100, min scanning window hit rate – 0.05, max missing calls

236    – 20.

237         Private allele counts per river were tallied using the populations module in

238    Stacks [31, 32] version 2.54 with default parameters. Populations with more than five

239    individuals were randomly subsampled to five to reduce the influence of uneven

240    sampling on the number of private alleles identified. Stacks was also used to calculate

241    other population level metrics such as observed heterozygosity, nucleotide diversity (Pi),

242    and Fis with default settings.

243         A genome-wide association (GWA) analysis was performed to characterize the

244    extent of isolation-by-distance previously reported by several authors (e.g., [7, 33]). The

245    trait of interest under investigation was latitude (the Chile strain of coho salmon was

246    excluded from this analysis). We used PLINK with default settings to perform this

247    analysis. Population structure was not included in this analysis as a covariate because

248    we were trying to characterize the fraction of the genome with a north – south gradient

249    and adding this covariate would remove much of that variation. R [34] and the qqman

250    package in R [35] were used to visualize the GWA analysis.

251         We tested for gene ontology enrichment based on the annotated variants that

252 were associated with the north – south gradient (for variants that were 'moderately'

253 likely to influence gene function and for those having 'low' or 'moderate' likelihood).

254 SnpEff [36] version 5.0e and the gene annotation from the NCBI were used to annotate

255 nucleotide variants for potential function alterations using default settings. Blast2GO

256 [37] and OmicsBox [38] version 2.0.36 were used to test for enriched GO categories

257 using default parameters.

258      To infer demographic histories of the salmon from the various rivers, we used

259 SMC++ [39] version 1.15.4.dev18+gca077da. In this analysis, we set the mutation rate to

260 8e-9 bp/generation and the generation time to 3 years. These parameters were

261 previously used in another coho salmon study examining demographic histories [7]. We

262 used nucleotide variants that were not filtered for rare variants (e.g., MAF < 0.05). We

263 also used the --missing-cutoff option (50 kbp) in SMC++ to reduce the influence of

264 missing genotypes (e.g., in centromeres).

265

266

267 **Results**

268 *Genome assemblies*

269      The size of both versions of the coho salmon genome assembly was 2.3 Gb,

270 which is also the same size of the closely related Chinook salmon (*O. tshawytscha*)

271 genome assembly [40]. However, version 2 of the coho salmon genome assembly was

272 much more contiguous than version 1 and had a more complete gene set (inferred from

273 BUSCO completeness). There was an almost 20x fold increase in contig N50 between

274 version 1 (58 kb) and version 2 (1,159 kb) of the genome assembly (Table 2). Likely as a

275 consequence of the increase in contiguity, the number of complete BUSCOs rose from

276 91% to 99%, which is comparable to the human genome assembly at 99% [41]. The

277 proportion of repeats also rose from 44.82% to 53.12% (compared to 52.94% in Chinook

278 salmon), and the number of annotated genes increased from 41,179 to 60,330 (47,105

279 in Chinook salmon). The NCBI reported that from version 1 to version 2, 37% of the

10

280    genome annotations were new and that 16% of the annotations on version 2 required

281    major changes from the previous version [42]. We note that the genome assembly was

282    produced from sequence data from two coho salmon and therefore not haplotype

283    resolved but chimeric in nature.

284         The coho salmon genome has extensive signatures of chromosomal duplication

285    (Figure 1, Table 2), which have been retained from the whole genome duplication

286    common to all salmonids [17]. The majority of duplicated regions from the salmonid-

287    specific genome duplication have diverged to a point where it is relatively easy to

288    differentiate between the copies (Figure 1, ≤90% identity), but certain sections of the

289    genome have retained high-sequence similarity where it is difficult to distinguish

290    between copies (Figure 1). Regions with very high-sequence similarity remain as

291    unplaced scaffolds as it was not possible to resolve which sequence belonged to which

292    duplicated region (see assembly methods; available on the NCBI website [43]). The

293    number of duplicate BUSCOs increased from 37% to 42.2% between versions (Table 2),

294    which suggests that the second assembly was able to distinguish between similar

295    paralogs/homeologs better whereas the first assembly likely collapsed them into a

296    single gene/BUSCO.

297         The coho salmon genome also has a high retention of repetitive elements (Figure

298    1, Table 2), which is another commonality of studied salmonids (e.g., [12, 18]). This is

299    especially true in regions near the centromere where the fraction of repetitive elements

300    is roughly 75% (Figure 1). That value is high compared to the genome average of 53%

301    (Table 2). For comparison, the most recent version of the Chinook salmon genome also

302    has a repeat content of 53% [44].

303

304    *Population genomics*

305         A PCA of 83 resequenced coho salmon genomes sampled from across North

306    America (and aquaculture samples), revealed that coho salmon clustered by region with

307    the exceptions of the Salmon River and Inch Creek (Figure 2). On the first principal

11

308    component of the PCA, the Salmon River clustered away from all the other samples. This

309    river belongs to the Thompson River watershed, and coho salmon from this region have

310    previously been observed to cluster in a similar manner [7]. Inch Creek salmon might

311    cluster separately as an artifact since the genome assembly was derived from an Inch

312    Creek salmon. This might increase read-alignment scores and influence SNP-calling in

313    some regions of the genome.

314        Excluding the Salmon River and Inch Creek samples, all other samples clustered

315    by region and by latitude in a manner consistent with isolation-by-distance suggested by

316    [7]. The Salmon River group have the lowest private allele counts (1,876 vs. a median of

317    4,188) and observed heterozygosity (0.22966 vs. a median of 0.285565). They also have

318    the highest total runs of homozygosity (Figures S1 and S2). The region with the highest

319    private allele count appears to be around the Puget Sound (e.g., Wallace River, private

320    allele count = 5,546) and Strait of Georgia regions (e.g., Capilano River, private allele

321    count = 6,415). Most of the northern rivers have low private allele counts with the

322    exception of the Kitimat River (private allele count = 6,341), which has the second

323    highest count (Figure 2).

324        To investigate how much of the genome has been influenced by isolation-by-

325    distance, we quantified the number of SNPs associated with the latitude gradient

326    observed from the PCA above (Figure 3, File S2). Roughly 13.9-33.8% of the 5,631,459

327    variants were associated with latitude at a significance level of 0.01-0.1 without multiple

328    test corrections (Figure 3). The proportion of variants associated with latitude dropped

329    to 0.07% after the alpha threshold was set to 0.05 with a Bonferroni correction (these

330    variants were widely distributed throughout the genome).

331        In Table 3, the most common nucleotide variant annotations from SNPeff are

332    shown, with intronic and intergenic variants being the most common type of variant

333    annotation. The variants that were significantly associated with latitude (see previous

334    paragraph, 0.07%) have a similar broad distribution of annotations relative to the entire

335    genome rather than enriched for variants that are likely to influence gene function

12

336   (Table 3, File S2). For instance, the percent of intergenic nucleotide variants remained at

337   31.2% of the total number of variants for the whole genome and for variants that were

338   significantly associated with latitude (Table 3). We would expect that if variants were

339   influencing traits under selection (e.g., based on latitude), the distribution would change

340   between all variants and those significantly associated with latitude if those SNPs

341   influenced gene function (e.g., 3' UTR and missense annotations).

342        Significant latitude-associated nucleotide variants (0.07%) identified in the GWA

343   analysis that were annotated as having a 'Moderate' likelihood to influence gene

344   function by SnpEff were found in 45 genes (File S2). Of these 45 genes, 4 genes had two

345   or more variants that were annotated as 'Moderate' in their impact on gene function

346   (Figure 4). No enriched gene ontologies were identified from genes with 'Moderate' or

347   even 'Moderate' + 'Low' (87 genes) nucleotide variant annotations (File S2). Only when

348   all genes with associated variants were tested, regardless of influencing function, did we

349   observe enriched GO terms (data not shown).

350        To put the nucleotide variation generated by isolation-by-distance into a broader

351   context, we identified possible times when northern populations could have recolonized

352   after the last glaciation period. By modeling demographic histories from genome

353   sequences using the SMC++ program, we were able to identify major decreases in

354   effective population size (Ne) that correspond with the Cordilleran Ice Sheet maximum

355   and the presumed penultimate global glacial maximum (Figure 5). We also observed

356   that for some populations, mostly northern, there was an additional drop in effective

357   population size between 3,750 and 8,000 years ago (Figure 5).

358

359

360   **Discussion**

361        As with previous analyses of salmonid genomes [10, 12, 17, 18, 45, 46], the

362   retention of duplicated chromosomes (i.e., homeologs) from the salmonid-specific

363   whole genome duplication [17] is a defining feature of the coho salmon genome. Some

13

364    of the duplicated regions have likely retained very high sequence similarity for roughly

365    90 million years (time estimate from [17, 45, 47]). A possible mechanism for high

366    sequence similarity retention is through tetrasomic inheritance [48].

367        The second version of the coho salmon genome assembly resolved a greater

368    number of duplicated regions of the genome compared to the first version. The better

369    resolution of duplicated regions can be observed with the increase in gene count and

370    the number of duplicated BUSCOs identified. Finer detail in these regions may help us in

371    future studies to better understand the residual impacts of whole genome duplication

372    on the biology of salmon.

373        From resequenced coho salmon genomes, we were able to better understand

374    population structure of coho salmon and its relationship with isolation-by-distance. One

375    of the striking features of the PCA of coho salmon populations was how divergent

376    Salmon River salmon were to all other populations. The Salmon River is part of the

377    Thompson River watershed and coho salmon from this system were thought to be

378    isolated from all other populations for potentially 150,000 years before secondary

379    contact roughly 13,500 years ago (essentially during the previous glacial period) [7]. This

380    would be consistent with findings in kokanee (*O. nerka*, a landlocked sockeye salmon

381    ecotype) in the upper Columbia River that similarly appear divergent from all other

382    populations of sockeye salmon and kokanee [12]. Taken together, these pieces of

383    evidence might be interpreted as support for a glacial refugium near the intersection of

384    the Cordilleran Ice Sheet and the Laurentide Ice Sheet.

385        A more likely alternative is that another unknown factor was influencing past

386    analyses and the PCA from the current study. The Salmon River coho salmon have

387    increased runs of homozygosity, reduced heterozygosity, and reduced private alleles,

388    which are indicators of a recent and extensive bottleneck. We were also able to infer

389    the demographic history from whole genome sequences of the Salmon River coho

390    salmon and found evidence of a bottleneck (from ~Ne 16,227 to ~Ne 1,749) around

391    4,000 years ago. These results could help explain why the Salmon River coho salmon

14

392    appear so divergent in a PCA, as low genetic diversity might be expected to increase the

393    amount of variation in the analysis since most of the other individuals do not have low

394    genetic diversity. We only collected samples from one tributary of the Thompson River

395    (a part of a much larger basin) and can only suggest that a plausible hypothesis from this

396    data is that recolonization of the Salmon River from a small founding population took

397    place after glaciers receded. We did not account for the influence of hatcheries, which

398    could also influence many of the metrics discussed above. Also, we did not incorporate

399    linked selection in demographic modeling as the type of analysis that we used was not

400    amenable. Without linked selection accounted for, there could be biases in times and

401    effective population sizes from our estimates [49]. Based on all the genetic diversity

402    metrics (above), demographic modeling, and what has previously been published on the

403    time of the most recent glacial maximum [50], however, recent recolonization of the

404    Salmon River remains a likely alternative hypothesis to a glacial refugium between ice

405    sheets.

406          Most other streams, except Inch Creek, clustered in the PCA based largely on

407    latitude for both PC1 and PC2 of the PCA. With a much more extensive sampling

408    strategy, Rougemont et al. (2020) found a similar trend and tested various demographic

409    histories [7]. The authors of that study found that the best supported model was a

410    glacial refugia to the south with recolonization of the northern streams after glacial

411    retreat – generating genomic signatures of isolation-by-distance. The private allele

412    analyses from the current study also supports this interpretation. The private allele

413    analysis identified that most of the northern streams have low private allele counts

414    compared to southern streams.

415          To better understand the pattern of genomic isolation-by-distance along the

416    latitude gradient, we performed a GWA analysis based on stream latitude. We found

417    that 13.9% of the variants were associated with latitude based on a $p$-value of 0.01

418    without multiple test correction. To put this into perspective, the north – south gradient

419    likely formed after the last Cordilleran Ice Sheet maximum 19-20,000 years ago [50] but

15

420    could have formed later between 3,750-8,000 years ago based on our demographic

421    history modeling. This would suggest that a large fraction of the nucleotide variants

422    responded within less than 8,000 years to the influence of isolation-by-distance.

423         We analyzed the influence of significantly associated nucleotide variants on gene

424    function and also the distribution of annotated variants to better understand if selection

425    played a large role in establishing the north – south gradient. We tested if genes with

426    significantly associated variants (α = 0.05, Bonferroni-correction), with 'Low' to

427    'Moderate' likelihoods of influencing gene function, belonged to any enriched GO

428    categories. If a trait was under selection based on latitude, we might expect enriched

429    GO terms associated with that trait. We did not find any enriched GO categories from

430    the GWA analysis. This may indicate that selection may not have contributed much to

431    establishing the north – south gradient.

432         When comparing the distribution of the most common variant annotations of

433    the full dataset with the variants that were significantly associated with latitude, the

434    largest fold difference was between the missense mutations (Table 3). We observed a

435    ~2x difference from 0.8% missense mutation rate in the entire genome to 1.7% in the

436    variants associated with latitude. While this increase does suggest selection may have

437    contributed to the development of the latitude gradient, we interpret that, because the

438    majority of the other annotations have similar frequencies, the majority of the variants

439    that make up the gradient are not under direct selection. Further, the increase in

440    missense mutations may represent slightly deleterious variants that escaped selective

441    pressure during postglacial recolonization and expansion. Linked selection could still

442    play a larger role but was not investigated here.

443         With or without selection, the north – south gradient of nucleotide variants likely

444    influences some phenotypic differences in a similar gradient. As an example, we

445    identified genes that had multiple nucleotide variants that are predicted to moderately

446    influence gene function, and which also have an association with latitude. These

447    included the rhotekin-like (*RTKN*, unknown function [51]), plectin-like (*PLEC*, giant

16

448    cytoskeleton scaffold [52]), PH and SEC7 domain-containing protein 4-like (*PSD4*, tight

449    junctions maintenance [53]), and GTPase IMAP family member 9-like (*Gimap9*, possibly

450    involved in T-cell development [54, 55]) genes. The nucleotide diversity of these genes

451    largely arises from the frequency of the alternative allele in the four most northern

452    streams – regions that would have likely been recolonized most recently assuming a

453    main southern glacial refugium.

454        Interestingly, *PLEC* is a candidate gene associated with migration distance in

455    brown trout (*Salmo trutta*), perhaps through its role in osmoregulation [56]. Cells

456    without *PLEC* were found to be more sensitive to changes in osmolarity (shrinking more

457    after exposure to urea) [57], and hatch-stage whitefish (*Coregonus lavaretus*) exposed

458    to high salinity have significantly higher *PLEC* protein expression [58]. Two of the other

459    genes with multiple variants moderately-likely to modify gene function and which were

460    associated with the north – south latitudinal gradient, *RTKN* and *PSD4*, may also have

461    roles in salinity tolerance. An Atlantic cod (*Gadus morhua* L.) nucleotide variant in the

462    intron of *PSD4* was found to be associated with a salinity gradient between the North

463    Sea and the Baltic Sea [59]. Likewise, researchers discovered that Rho (RTKN is an

464    effector protein of RhoA [51, 60]), is activated by hyperosmotic stress [61].

465        *PLEC*, *PSD4*, and *RTKN* all appear to be involved in cell junction functionality. Cell

466    junctions observed in a *PLEC* knockout cell line appeared to be compromised [62], *PSD4*

467    (also known as EFA6B) is required for efficient tight junction formation [63], and *RTKN*

468    influences cell junctions through PIST [51] and Septin proteins [51, 64]. It is thought that

469    cellular tight junctions play an important role in water and salt balance in teleost fishes

470    [65]. Considering that three of the four genes with multiple latitude associated

471    nucleotide variants and which are moderately likely to alter gene function could impact

472    cell junctions and osmoregulation, we hypothesize that ocean salinity may have

473    influenced coho salmon recolonization of northern streams.

474        While, the Pacific Ocean salinity was thought to be ~4% higher during the last

475    glacial maximum as freshwater was stored in glaciers [66], it is difficult to predict how

17

476  the salinity gradient observed in modern times [67] might have been influenced as

477  glaciers began to retreat (when northern recolonization would have been possible). If

478  there was a difference in salinity between northern and southern regions, nucleotide

479  variation in these genes may have facilitated northern colonization in some way.

480       From inferred demographic histories, we were able to estimate a recolonization

481  date of some northern streams (based on the founder effect that would be expected to

482  accompany recolonization) to between 3,750-8,000 years ago. This places an upper limit

483  on the age of the latitude gradient and how swiftly such a gradient can form. These

484  values are based on assumptions of a mutation rate of 8e-9 bp/generation and a

485  generation time of three years. Linked selection may also bias our time and effective

486  population estimates [49] as we did not account for them in modeling.

487       While it is important to remember that time and population estimates are

488  influenced by many factors when inferring demographic histories from sequence data,

489  multiple lines of evidence can be used to strengthen these inferences or put them in a

490  more realistic context. Radiometric evidence supports that the Cordilleran Ice Sheet

491  maximum occurred between 19,000 and 20,000 years ago [50]. Likewise, chemical

492  properties of gases in Antarctic ice cores support this termination of the last glaciation

493  period (Termination I) to roughly the same time period, as well as a previous

494  termination of the penultimate glaciation period around 138,000 and 148,000 years ago

495  (Termination II) [68]. In the demographic histories of the coho salmon, we noted

496  dramatic declines of nearly all salmon populations for both these time periods. This

497  observation supports the parameters used for modeling the demographic histories as

498  we expect that populations might decline in response to increased glaciation or rapid

499  climate change.

500       The overall trend we observed from modeling demographic histories was major

501  drops in effective population size at each transition from glaciation to inter-glaciation

502  period with increases for nearly all populations after the penultimate glaciation period

503  and uncommon increases for specific rivers after the most recent glaciation period. At a

18

504 species level, these transitional drops likely influence multiple aspects of coho salmon

505 biology since genetic variability can contribute to many characteristics of a species.

506

507 **Conclusions**

508       In this study, we generated two reference genome assemblies as tools for

509 conservation and management of coho salmon. Additionally, we resequenced the

510 genomes of a wide distribution of coho salmon from rivers along North America. We

511 were able to identify a north – south gradient in the nucleotide variation of the

512 genomes, which had been observed in previous studies. To add to previous

513 observations, we quantified that approximately 13.9% of the variation in the

514 resequenced genomes followed the north – south gradient. We also were able to

515 estimate that the age of the north – south gradient is likely under 8,000 years of age.

516 This gradient likely contributes to phenotypic diversity between northern and southern

517 rivers since we identified gene modifying variants that were associated with the latitude

518 gradient. Finally, we modeled demographic histories of the coho salmon from different

519 rivers and discovered that major drops in effective population size were related to

520 changes between glacial and inter-glacial periods. We believe the coho salmon genome

521 assemblies will facilitate research to better understand coho salmon biology and may

522 enhance management of this culturally and economically important species.

523

524 **Data availability**

525 Raw data for the genome assembly was submitted to the NCBI under the BioProject

526 PRJNA352719. Whole genome resequencing data was submitted under PRJNA401427

527 and PRJNA808051 to the NCBI BioProjects (see File S1 for specific samples used in this

528 study). The VCF file used for analyses in this study was submitted to the GSA Figshare

529 portal.

530

531 **Acknowledgements**

      19

20

560     provided to support production of doubled haploids.

561

562     **Author Contributions**

563         EBR, KAC, JSL, and BFK performed genome assembly, chromosome assembly,

564     genome submission, and generated genome metrics. DRM performed repeat library

565     construction. EBR, CAD, MTC, AM, and DS performed wet-lab work including DNA and

566     RNA extractions and mitochondrial sequencing. EBR, QR, EN, DRM, and JSL performed

567     SNP calling and population genomic analyses. RHD, MTC, DS, and CB generated, raised,

568     and dissected doubled-haploid samples for the genome assembly and transcriptome.

569     REB, TDB, KAN, and JMY provided samples used in resequencing work. KAN provided

570     early access to linkage map and additional guidance on its use. RHD, REW, TDB, KAN,

571     JMY, RN, LB, WSD, SJMJ, and BFK initiated, planned and supervised the project. EBR,

572     KAC, DRM and BFK wrote first draft of the manuscript.

573

574

575     **References**

1. Irvine JR, Fukuwaka M. Pacific salmon abundance trends and climate change. ICES J Mar Sci. 2011;68:1122–30.

2. Beamish RJ, Sweeting RM, Lange KL, Noakes DJ, Preikshot D, Neville CM. Early Marine Survival of Coho Salmon in the Strait of Georgia Declines to Very Low Levels. Mar Coast Fish. 2010;2:424–39.

3. Beamish RJ, McFarlane GA, Thomson RE. Recent declines in the recreational catch of coho salmon (Oncorhynchus kisutch) in the Strait of Georgia are related to climate. Can J Fish Aquat Sci. 2011. https://doi.org/10.1139/f98-195.

4. Kristianson G, Strongitharm D. The Evolution of Recreational Salmon Fisheries in British Columbia. Vancouver, BC: Pacific Fisheries Resource Conservation Council; 2006.

5. Noakes DJ, Beamish RJ, Sweeting R, King J. Changing the balance: Interactions between hatchery and wild Pacific coho salmon in the presence of regime shifts. North Pac Anadromous Fish Commision Bull. 2000;:155–63.

6. Beacham TD, Wetklo M, Deng L, MacConnachie C. Coho Salmon Population Structure in North America Determined from Microsatellites. Trans Am Fish Soc. 2011;140:253–70.

7. Rougemont Q, Moore J-S, Leroy T, Normandeau E, Rondeau EB, Withler RE, et al. Demographic history shaped geographical patterns of deleterious mutation load in a broadly distributed Pacific Salmon. PLOS Genet. 2020;16:e1008348.

8. Rougemont Q, Xuereb A, Dallaire X, Moore J-S, Normandeau E, Rondeau EB, et al. Long-distance migration is a major factor driving local adaptation at continental scale in Coho salmon. Mol Ecol. 2022;n/a n/a.

9. Zhang H, Tan E, Suzuki Y, Hirose Y, Kinoshita S, Okano H, et al. Dramatic improvement in genome assembly achieved using doubled-haploid genomes. Sci Rep. 2014;4:6780.

10. Christensen KA, Leong JS, Sakhrani D, Biagi CA, Minkley DR, Withler RE, et al. Chinook salmon (Oncorhynchus tshawytscha) genome and transcriptome. PLOS ONE. 2018;13:e0195461.

11. Rondeau EB, Minkley DR, Leong JS, Messmer AM, Jantzen JR, Schalburg KR von, et al. The Genome and Linkage Map of the Northern Pike (Esox lucius): Conserved Synteny Revealed between the Salmonid Sister Group and the Neoteleostei. PLOS ONE. 2014;9:e102089.

12. Christensen KA, Rondeau EB, Minkley DR, Sakhrani D, Biagi CA, Flores A-M, et al. The sockeye salmon genome, transcriptome, and analyses identifying population defining regions of the genome. PLOS ONE. 2020;15:e0240935.

13. Gnerre S, MacCallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, et al. High-quality draft assemblies of mammalian genomes from massively parallel sequence data. Proc Natl Acad Sci. 2011;108:1513–8.

14. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the Gap: Upgrading Genomes with Pacific Biosciences RS Long-Read Sequencing Technology. PLOS ONE. 2012;7:e47768.

15. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. Bioinformatics. 2015;31:3210–2.

16. Krzywinski MI, Schein JE, Birol I, Connors J, Gascoyne R, Horsman D, et al. Circos: An information aesthetic for comparative genomics. Genome Res. 2009. https://doi.org/10.1101/gr.092759.109.

17. Allendorf FW, Thorgaard GH. Tetraploidy and the Evolution of Salmonid Fishes. In: Turner BJ, editor. Evolutionary Genetics of Fishes. Springer US; 1984. p. 1–53.

18. Lien S, Koop BF, Sandve SR, Miller JR, Kent MP, Nome T, et al. The Atlantic salmon genome provides insights into rediploidization. Nature. 2016;533:200–5.

19. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. Repbase Update, a database of eukaryotic repetitive elements. Cytogenet Genome Res. 2005;110:462–7.

20. Smit A, Hubley R. RepeatModeler Open-1.0. 2013. http://www.repeatmasker.org/. Accessed 18 Dec 2017.

21. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. BMC Bioinformatics. 2009;10:1–9.

22. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. ArXiv13033997 Q-Bio. 2013.

23. github.com/broadinstitute/picard. Java. Broad Institute; 2020.

24. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A, et al. From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. Curr Protoc Bioinforma. 2013;43:11.10.1-11.10.33.

25. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. Genome Res. 2010;20:1297–303.

26. Barría A, Christensen KA, Yoshida G, Jedlicki A, Leong JS, Rondeau EB, et al. Whole Genome Linkage Disequilibrium and Effective Population Size in a Coho Salmon (Oncorhynchus kisutch) Breeding Population Using a High-Density SNP Array. Front Genet. 2019;10:498.

27. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and VCFtools. Bioinformatics. 2011;27:2156–8.

28. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve years of SAMtools and BCFtools. GigaScience. 2021;10:giab008.

29. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to the challenge of larger and richer datasets. GigaScience. 2015;4.

30. PLINK 1.9. http://www.cog-genomics.org/plink/1.9/. Accessed 1 Jun 2018.

31. Catchen J, Hohenlohe PA, Bassham S, Amores A, Cresko WA. Stacks: an analysis tool set for population genomics. Mol Ecol. 2013;22:3124–40.

32. Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH. Stacks: Building and Genotyping Loci De Novo From Short-Read Sequences. G3 GenesGenomesGenetics. 2011;1:171–82.

33. Smith CT, Nelson RJ, Wood CC, Koop BF. Glacial biogeography of North American coho salmon (Oncorhynchus kisutch). Mol Ecol. 2001;10:2775–85.

34. R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria: R Foundation for Statistical Computing; 2020.

35. Turner SD. qqman: an R package for visualizing GWAS results using Q-Q and manhattan plots. bioRxiv. 2014;:005165.

36. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster strain w1118; iso-2; iso-3. Fly (Austin). 2012;6:80–92.

37. Götz S, García-Gómez JM, Terol J, Williams TD, Nagaraj SH, Nueda MJ, et al. High-throughput functional annotation and data mining with the Blast2GO suite. Nucleic Acids Res. 2008;36:3420–35.

38. OmicsBox - BioBam | Bioinformatics Made Easy. BioBam. https://www.biobam.com/omicsbox/. Accessed 1 Mar 2022.

39. Terhorst J, Kamm JA, Song YS. Robust and scalable inference of population history from hundreds of unphased whole genomes. Nat Genet. 2017;49:303–9.

40. Oncorhynchus tshawytscha genome assembly Otsh_v2.0. NCBI. https://ncbi.nlm.nih.gov/data-hub/assembly/GCF_018296145.1/. Accessed 6 Apr 2022.

41. Homo sapiens genome assembly GRCh38.p13. NCBI. https://ncbi.nlm.nih.gov/data-hub/assembly/GCF_000001405.39/. Accessed 6 Apr 2022.

42. Oncorhynchus kisutch Annotation Report. https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Oncorhynchus_kisutch/101/. Accessed 30 Mar 2022.

43. Oncorhynchus kisutch genome assembly Okis_V2. NCBI. https://ncbi.nlm.nih.gov/data-hub/assembly/GCF_002021735.2/. Accessed 6 Apr 2022.

44. Oncorhynchus tshawytscha Annotation Report. https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Oncorhynchus_tshawytscha/101/. Accessed 6 Apr 2022.

45. Berthelot C, Brunet F, Chalopin D, Juanchich A, Bernard M, Noël B, et al. The rainbow trout genome provides novel insights into evolution after whole-genome duplication in vertebrates. Nat Commun. 2014;5.

46. De-Kayne R, Zoller S, Feulner PGD. A de novo chromosome-level genome assembly of Coregonus sp. "Balchen": One representative of the Swiss Alpine whitefish radiation. Mol Ecol Resour. 2020;20:1093–109.

47. Macqueen DJ, Johnston IA. A well-constrained estimate for the timing of the salmonid whole genome duplication reveals major decoupling from species diversification. Proc R Soc B Biol Sci. 2014;281.

48. Allendorf FW, Bassham S, Cresko WA, Limborg MT, Seeb LW, Seeb JE. Effects of Crossovers Between Homeologs on Inheritance and Population Genomics in Polyploid-Derived Salmonid Fishes. J Hered. 2015;106:217–27.

49. Pouyet F, Aeschbacher S, Thiéry A, Excoffier L. Background selection and biased gene conversion affect more than 95% of the human genome and bias demographic inferences. eLife. 2018;7:e36317.

50. Clark PU, Dyke AS, Shakun JD, Carlson AE, Clark J, Wohlfarth B, et al. The Last Glacial Maximum. Science. 2009;325:710–4.

51. Ito H, Morishita R, Nagata K. Functions of Rhotekin, an Effector of Rho GTPase, and Its Binding Partners in Mammals. Int J Mol Sci. 2018;19:2121.

52. Winter L, Wiche G. The many faces of plectin and plectinopathies: pathology and mechanisms. Acta Neuropathol (Berl). 2013;125:77–93.

53. Zangari J, Partisani M, Bertucci F, Milanini J, Bidaut G, Berruyer-Pouyet C, et al. EFA6B Antagonizes Breast Cancer. Cancer Res. 2014;74:5493–506.

54. Wang Z, Li X. IAN/GIMAPs are conserved and novel regulators in vertebrates and angiosperm plants. Plant Signal Behav. 2009;4:165–7.

55. Limoges M-A, Cloutier M, Nandi M, Ilangumaran S, Ramanathan S. The GIMAP Family Proteins: An Incomplete Puzzle. Front Immunol. 2021;12.

56. Lemopoulos A, Uusi-Heikkilä S, Hyvärinen P, Alioravainen N, Prokkola JM, Elvidge CK, et al. Association Mapping Based on a Common-Garden Migration Experiment Reveals

25

Candidate Genes for Migration Tendency in Brown Trout. G3 GenesGenomesGenetics. 2019;9:2887–96.

57. Osmanagic-Myers S, Gregor M, Walko G, Burgstaller G, Reipert S, Wiche G. Plectin-controlled keratin cytoarchitecture affects MAP kinases involved in cellular stress response and migration. J Cell Biol. 2006;174:557–68.

58. Papakostas S, Vasemägi A, Vähä J-P, Himberg M, Peil L, Primmer CR. A proteomics approach reveals divergent molecular responses to salinity in populations of European whitefish (Coregonus lavaretus). Mol Ecol. 2012;21:3516–30.

59. Berg PR, Jentoft S, Star B, Ring KH, Knutsen H, Lien S, et al. Adaptation to Low Salinity Promotes Genomic Divergence in Atlantic Cod ( Gadus morhua L.). Genome Biol Evol. 2015;7:1644–63.

60. Thumkeo D, Watanabe S, Narumiya S. Physiological roles of Rho and Rho effectors in mammals. Eur J Cell Biol. 2013;92:303–15.

61. Ciano-Oliveira CD, Sirokmány G, Szászi K, Arthur WT, Masszi A, Peterson M, et al. Hyperosmotic stress activates Rho: differential involvement in Rho  kinase-dependent MLC phosphorylation and NKCC activation. Am J Physiol-Cell Physiol. 2003;285:C555–66.

62. Jirouskova M, Nepomucka K, Oyman-Eyrilmez G, Kalendova A, Havelkova H, Sarnova L, et al. Plectin controls biliary tree architecture and stability in cholestasis. J Hepatol. 2018;68:1006–17.

63. Théard D, Labarrade F, Partisani M, Milanini J, Sakagami H, Fon EA, et al. USP9x-mediated deubiquitination of EFA6 regulates de novo tight junction assembly. EMBO J. 2010;29:1499–509.

64. Kim J, Cooper JA. Junctional Localization of Septin 2 Is Required for Organization of Junctional Proteins in Static Endothelial Monolayers. Arterioscler Thromb Vasc Biol. 2021;41:346–59.

65. Chasiotis H, Kolosov D, Bui P, Kelly SP. Tight junctions, tight junction proteins and paracellular permeability across the gill epithelium of fishes: A review. Respir Physiol Neurobiol. 2012;184:269–81.

66. Insua TL, Spivack AJ, Graham D, D'Hondt S, Moran K. Reconstruction of Pacific Ocean bottom water salinity during the Last Glacial Maximum. Geophys Res Lett. 2014;41:2914–20.

67. Durack PJ, Wijffels SE, Matear RJ. Ocean Salinities Reveal Strong Global Water Cycle

26

Intensification During 1950 to 2000. Science. 2012;336:455–8.

68. Kawamura K, Parrenin F, Lisiecki L, Uemura R, Vimeux F, Severinghaus JP, et al. Northern Hemisphere forcing of climatic cycles in Antarctica over the past 360,000 years. Nature. 2007;448:912–6.

69. Wickham H. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York; 2016.

70. Becker RA, Wilks AR, Brownrigg R, Minka TP, Deckmyn A. maps: Draw Geographical Maps. 2018.

71. Chappellaz J, Barnola JM, Raynaud D, Korotkevich YS, Lorius C. Ice-core record of atmospheric methane over the past 160,000 years. Nature. 1990;345:127–31.

72. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. Genome Res. 2009;19:1655–64.

576

577 **Figures**

578 **Figure 1. Circos plot of the first version of the coho salmon genome assembly.** In the

579 interior of the Circos plot are the links between duplicated regions of the

580 chromosomes/linkage groups (i.e., homeologous regions). A) Representations of the

581 chromosomes with the approximate position of the centromere marked by a filled

582 circle. The tick marks represent 10 Mbp intervals. B) The percent identity between

583 duplicated regions of the chromosome. The orange-red color represents very high

584 similarity (> 94%), the orange color high similarity (91-94%), the yellow moderate (88-

585 91%), and the green low (< 88%). C) The fraction of repetitive elements, with red

586 representing high (> 60%), yellow as moderate (35-60%), and green as low (< 35%).

587

588 **Figure 2. Influence of isolation-by-distance on coho salmon population structure.** A) A

589 PCA of coho salmon based on variants that were filtered for linkage disequilibrium

590 (plotted using ggplot [69]). Circles were drawn around Inch Creek and Salmon River

591 individuals to highlight that much of the variation of this PCA was due to differences of

592  salmon from these rivers. B) The same figure as A, with Salmon and Inch Creek salmon

593  removed. When Inch Creek and Salmon River salmon are removed from the graph, the

594  influence of latitude/isolation-by-distance can be observed on PC1 and PC2. Individuals

595  from the same river (same colour) also tended to cluster near each other. C) A map of

596  the various rivers sampled for this study and the corresponding private allele counts

597  (plotted with the maps [70] package in R).  The private allele counts are also displayed to

598  the side as a bar graph.

599

600  **Figure 3. Fraction of the genome responsive to isolation-by-distance.** A) A Manhattan

601  plot of the latitude genome-wide association analysis of all coho salmon except for the

602  Chile strain. The red line represents a significance level of 0.01 after a Bonferroni

603  correction and the blue line a 0.05 level with the same correction. B) A histogram of

604  variant counts for different significance levels. C) The percent of the variants at different

605  significance levels. The percentage represents all the variants at or above the line. There

606  are 33.8% of the variants that have a significant association with latitude at $p \leq 0.1$

607  without multiple test correction.

608

609  **Figure 4. Genes with multiple 'moderate' nucleotide variants that may influence**

610  **function.** Four genes were found to be associated with latitude and also contained

611  multiple SNPs that likely modify gene function. The four genes are: rhotekin-like (RTKN

612  LOC109895613), plectin-like (PLEC LOC109904478), PH and SEC7 domain-containing

613  protein 4-like (PSD4 LOC109868337), and GTPase IMAP family member 9-like (Gimap9

614  LOC109880231). A) Pie diagrams showing the distribution of reference (blue) and

615  alternative alleles (red) for each gene and location. B) Map produced with the maps

616  package in R showing the sampling sites.

617

618  **Figure 5. Demographic histories of coho salmon populations based on genome**

619  **sequencing.** A) Each labeled line represents multiple individuals from the same river or

28

620    strain. The X-axis represents calendar years based on a generation time of 3 years for

621    coho salmon.  The Y-axis is the effective population size (Ne) estimate. The estimated

622    age of the Cordilleran ice sheet maximum was taken from [50]. The approximate age of

623    the last interglacial period was based on [71]. B) For each location, a number nearby

624    indicates a drop of at least 5,000 in the Ne for one of the time points noted in A. The

625    colour of the river label indicates if the river had a drop in Ne during the 1$^{st}$ time interval

626    (orange – yes, blue – no, the Klamath River was the only southern river with a drop

627    during the 1$^{st}$ time period).

628

629    **Tables**

630    **Table 1.** Whole-genome resequencing sources

| Source | Country | State/Province | Count |
| --- | --- | --- | --- |
| | | | Female, Male |
| Klamath River (Hatchery) | US | CA/OR | 1F, 4M |
| Deschutes River (Hatchery) | US | CA/OR | 2F, 3M |
| Big Quilcene River (Hatchery) | US | WA | 2F, 3M |
| Wallace River (Hatchery) | US | WA | 6M, 4? |
| Tsoo-Yess River (Hatchery) | US | WA | 1F, 4M |
| Inch Creek (Hatchery) | Canada | BC | 3F, 5M |
| Capilano River (Hatchery) | Canada | BC | 5F |
| Robertson Creek (Hatchery) | Canada | BC | 5M |
| Salmon River (Hatchery) | Canada | BC | 5F |
| Pallant Creek | Canada | BC | 5M |
| Kitimat River (Hatchery) | Canada | BC | 5F, 5M |
| Berners River | US | AK | 2F, 3M |
| Kwethluk River | US | AK | 1F, 4M |
| AquaChile (Strain) | Chile | NA | 5F |

631    State/Province Abbreviations: CA - California, OR - Oregon, WA - Washington, BC -

29

632    British Columbia, AK – Alaska

633

634    **Table 2.** Genome statistics

|  | Contig N50 | Contig # | BUSCO | % Repeats | Genes |
|---|---|---|---|---|---|
| Ver 1 | 58,118 | 97,074 | 91%-55:37* | 44.82† | 41,179† |
| Ver 2 | 1,159,298 | 8,770 | 99.2%-57.1:42.2*† | 53.12† | 60,330† |

635    Ver 1, NCBI: GCF_002021735.1; Ver 2, NCBI: GCF_002021735.2

636    *Percent complete-single:duplicate

†Reported by NCBI (NCBI used actinopterygii_odb10 for BUSCO)

637

638    **Table 3.** Distribution of common nucleotide variant annotations

|  | Entire genome* | Associated with latitude* |
|---|---|---|
| **Intron** | 44.0% | 42.1% |
| **Intergenic** | 31.2% | 31.2% |
| **Upstream** | 10.4% | 11.2% |
| **Downstream** | 7.3% | 7.6% |
| **3' UTR** | 2.8% | 3.2% |
| **Missense** | 0.8% | 1.7% |
| **Synonymous** | 1.1% | 1.1% |

639    *5,631,459 genomic SNPs, 3,940 significant SNPs from GWA analysis

640

641    **Supplemental Material**

642    **Figure S1. Runs of homozygosity and admixture among coho salmon from different**

643    **streams.** A) The top figure shows the total runs of homozygosity (ROH) for each

644    individual (default settings in PLINK). The bottom figures show the admixture of each

645    individual based on cluster counts of k=2 and k=3 (with Admixture software [72] using

30

646    default settings with LD filtered SNPs). Streams are shown at the bottom and delineated

647    by the alternating blue bar. B) A map (generated using the maps package in R) showing

648    the locations in A.

649

650    **Figure S2. The relationship between runs of homozygosity and latitude.** A) Counts of

651    runs of homozygosity (ROH) and the total length of the ROH when combined (see

652    Methods for parameters used). There is a distinct cluster of Salmon River individuals

653    with higher counts and lengths of ROH. B) The relationship between the average length

654    of ROH per individual and latitude. The line was plotted using the geom_smooth

655    function in ggplot2 with the linear model method. Latitude significantly ($p = 0.029$)

656    explained variation in the average length of ROH (~6% of the variation, Adjusted $R^2$ =

657    0.04886).

658

659    **File S1. Sample information and SRA accession numbers.**

660

661    **File S2. Nucleotide variants significantly associated with latitude.** The

662    SignificantVariants tab in this spreadsheet file has information on all of the significantly

663    associated SNPs with latitude. The ModerateGenes tab has information on SNPs that

664    were both associated with latitude and also have annotations from SNPeff that were

665    moderately likely to influence gene function. The Moderate+LowGenes tab has

666    information on SNPs that were both associated with latitude and also have an

667    annotation from SNPeff that were likely to have moderate or low influences on gene

668    function. The GO tab has two lists of genes that were used in the GO enrichment

669    analyses. The Frequency of MultiVariantGenes tab has information on the genes with

670    multiple SNPs thought to moderately influence gene function and which were also

671    associated with latitude. This information was used to generate pie charts. The

672    GenotypeGenes tab has the genotypes for each individual for genes in the Frequency of

673    MultiVariantGenes tab. The Regions tab has information on the latitudes used for each

674     stream. The DistributionOfVariantAnnotations tab has information on SNPeff

675     annotations from the SNPs that were significantly associated with latitude as well as the

676     SNPs from the entire genome.