1  **Genome assembly, transcriptome and SNP database for chum salmon**

2  **(*Oncorhynchus keta*)**

3

4  **Eric B. Rondeau[1,2,3*], Kris A. Christensen[1,2*], Dionne Sakhrani[1], Carlo A. Biagi[1], Mike**

5  **Wetklo[3], Hollie A. Johnson[2], Cody A. Despins[2], Rosalind A. Leggatt[1], David R. Minkley[2],**

6  **Ruth E. Withler[3], Terry D. Beacham[3], Ben F. Koop[2], Robert H. Devlin[1§]**

7

8

9  [1]Fisheries and Oceans Canada, 4160 Marine Dr., West Vancouver, British Columbia, V7V 1N6,

10  Canada

11  [2]Department of Biology, University of Victoria, Victoria, British Columbia, V8W 3N5, Canada

12  [3]Pacific Biological Station, Fisheries and Oceans Canada, Nanaimo, British Columbia, V9T 6N7,

13  Canada

14

15  [§]Corresponding author

16  *Authors contributed equally to results of manuscript

17

18  Email addresses:

19  EBR: eric.rondeau@dfo-mpo.gc.ca

20  KAC: kris.christensen@wsu.edu

21  DS: Dionne.Sakhrani@dfo-mpo.gc.ca

22  CAB: Carlo.Biagi@dfo-mpo.gc.ca

23  MW: mike.wetklo@dfo-mpo.gc.ca

24  HAJ: holliej@uvic.ca

25  CAD: cdespins@uvic.ca

26  RAL: Rosalind.Leggatt@dfo-mpo.gc.ca

27  DRM: dminkley@uvic.ca

28  REW: RWithler@shaw.ca

29  TDB: Terry.Beacham@dfo-mpo.gc.ca

30  BFK: bkoop@uvic.ca

31    RHD: Robert.Devlin@dfo-mpo.gc.ca

**Abstract**

Chum salmon (*Oncorhynchus keta*) is the species with the widest geographic range of the anadromous Pacific salmonids,. Chum salmon is the second largest of the Pacific salmon, behind Chinook salmon, and considered the most plentiful Pacific salmon by overall biomass. This species is of significant commercial and economic importance: on average the commercial chum salmon fishery has the second highest processed value of the Pacific salmon within British Columbia. The aim of this work was to establish genomic baseline resources for this species. Our first step to accomplish this goal was to generate a chum salmon reference genome assembly from a doubled-haploid chum salmon. Gene annotation of this genome was facilitated by an extensive RNA-seq database we were able to create from multiple tissues. Range-wide resequencing of chum salmon genomes allowed us to categorize genome-wide geographic variation, which in turn reinforced the idea that genetic differentiation was best described on a regional, rather than at a stock-specific, level. Within British Columbia, chum salmon regional groupings were described at the conservation unit (CU) level, and there may be substructure within particular CUs. Genome wide associations of phenotypic sex to SNP genetic markers identified two clear peaks, a very strong peak on Linkage Group 15, and another on Linkage Group 3. With these new resources, we were better able to characterize the sex-determining region and gain further insights into sex determination in chum salmon and the general biology of this species.

**Keywords**

Chum salmon, Genome assembly, *Oncorhynchus keta*, RNA-seq, Resequencing, SNP database.

**Background**

Pacific salmon of the genus *Oncorhynchus* are iconic, culturally important keystone species spawning across freshwater watersheds that feed the Northern Pacific Ocean. Predominately anadromous, members of most species spend years at sea, consuming marine nutrients that are eventually deposited into coastal ecosystems where they provide a valuable source of food to numerous marine and terrestrial species as the salmon spawn and then die [1].

Chum salmon (*Oncorhynchus keta*) are the second largest of the Pacific salmonids and may have historically represented up to 50% of the salmonid biomass in the Pacific Ocean [2]. It is the most widely distributed of the Pacific salmonid species [3, 4], with spawning grounds ranging from Japan and the eastern coast of the Korean Peninsula through to Northern Russia, and from the Mackenzie River south through Central California in North America [5].   Among the most significant species of Pacific salmon in commercial fisheries – in an analysis of British Columbian commercial fisheries 2012-2015, chum salmon was the most plentiful species by weight in 3 out of 4 years analyzed, and second most valuable by processed value when averaged across the four year period ($31 million per year) [6].

A key and fascinating biological feature in salmonids is homing, whereby adults demonstrate an ability to return to the same riverine sites where they were spawned, although not all species show the same degree of site fidelity (reviewed in [7]). Some species, such as Sockeye, have been observed to return to within metres of where they were hatched (e.g., [8]), but other species vary in their fidelity to site of return and stray rate. Reasons for straying are likely varied (reviewed in [7]), but significant factors are thought to be juvenile freshwater residence time and freshwater migration distance, both of which lead to reduced imprinting. With chum salmon having relatively short freshwater residence (they migrate to sea as fry) and short migration distances (on average), it is perhaps not surprising that chum tend to have higher than average stray rates among the Pacific salmonids [7]. The consequences of such straying are that while regional-level differentiation (e.g., [9, 10]) and run-timing differentiation between summer and fall runs (e.g., [11–13]) can be observed, population-level genetic differentiation is not often seen within chum salmon.

86      The genomes of salmonids, including chum salmon, possess a key feature shared by all

87      salmonid genomes, a salmon-lineage specific whole-genome duplications (WGD). WGDs very

88      likely play one of the more significant roles in evolutionary innovation [14–17] and are found in

89      plants (reviewed in [18] ), fungi [19, 20] , arthropods [21, 22], basal vertebrates ~500 million

90      years ago (mya) [15, 23, 24], fishes ~300 mya [25–27], and more recently in ancestral salmonids

91      ~90 mya [28, 29] . These major genome expansions have been proposed to allow for

92      adaptations to new niches or conditions, particularly in times of major environmental change

93      (reviewed in [30]). The occurrence of over 70 different salmonid species lineages stemming

94      from the relatively recent ancestral WGD [29] offers a valuable system to i) observe

95      evolutionary consequences of a relatively recent autopolyploid WGD, ii) identify ensuing

96      mechanisms for regaining stable meiosis and cell division by regaining a functional diploid state

97      through re-diploidization, and iii) draw associations between mechanisms of re-diploidization to

98      potential genetic specialization that allow for species adaptation such as disease resistance.

99      Additionally, each species has evolved unique morphology, life history strategies, and responses

100     to common salmon pathogens (e.g., varied resistance to salmon aquaculture from pathogens

101     such as the sea louse [31, 32]). This phenotypic variety provides future opportunities for

102     exploring the biology and genetics behind the genomic architecture of whole-genome

103     duplication have shaped these unique species.

104     The presence of these duplications, however, can present major technological

105     challenges to genome assembly, due to limited differentiation between duplicated portions of

106     the genome. Salmonids offer additional hurdles in that a significant portion of the genome still

107     remains in a tetraploid-like state [33, 34], and may show lineage-specific re-diploidization

108     patterns [35], or chromosome architecture through species-specific fusions [36]. While many

109     challenges remain, the technological barriers to assembly of salmonid genomes are beginning

110     to fall, as evidenced by the relatively rapid recent release of salmonid genomes [37–44].  A

111     fully-annotated reference chum salmon genome will enhance development of genomics-based

112     technologies to improve the effectiveness of fisheries management of the wild chum salmon

113     fishery. This has already been performed for other Pacific salmon species in British Columbia

114 (e,g., [45, 46]), and a genome assembly for chum would provide the ability to adopt similar

115 management tools based on emerging high-throughput sequencing technologies.

116      Genetic resources in chum salmon have, as in many other species, been in a state of

117 transition as genetic tools have advanced and become more widespread. Early work on

118 population genetic structure in chum salmon utilized allozymes [47, 48] and microsatellite

119 markers [9, 49] and provided the first range-wide studies on genetic diversity [10]. Recently,

120 genetic stock identification tools have been shifting from microsatellites to single-nucleotide

121 polymorphisms (SNPs), providing increased accuracy of genetic discrimination with increasing

122 marker numbers [50]. Early identification of SNPs in chum salmon [51–55] led to the

123 development of a SNP panel for assessing genetic diversity and population structures in chum

124 salmon [13]; development of expanded SNP panels for fisheries management continues to

125 occur with increased marker density and improving genetic baselines allow for increased power

126 ([56; Beacham T.D. and Sutherland B.J.G, Personal Communication). Restriction-site Associated

127 DNA sequencing (RADseq) has recently enabled a much more rapid throughput for SNP

128 discovery [57, 58], and studies in chum have utilized this technique to enable researchers to

129 develop linkage maps to explore regions of residual inheritance associated with the

130 aforementioned genome duplication event [34]. This advance in technology has further allowed

131 for the identification of extended patterns of linkage disequilibrium, demonstrating the power

132 of increased marker density on the identification of genomic features of large effect [59].

133 Despite this significant effort, unlike in other *Oncorhynchus* species (e.g., Rainbow trout [60];

134 Chinook salmon [41], sockeye salmon [61]), neither a whole-genome catalog of SNP markers

135 nor whole-genome resequencing data has been available as a resource for chum salmon to

136 date. The development of such a resource will further allow genetic resources, such as SNP

137 panels, to be placed in context relative to genes or other annotated genomic features.

138      In this work, we have sequenced and assembled the genome of a mitotic gynogen

139 doubled haploid chum salmon to eliminate allelic variation but retain paralog differences.

140 Extensive multi-tissue RNA-seq was generated to provide the base for annotation of the

141 genome as well as a tissue-specific expression atlas for future comparative studies. Finally,

142 whole-genome resequencing was performed across 59 individual chum salmon from a select

6

143   distribution of the species' range to catalogue genome-wide diversity in this species. The utility

144   of the dataset is further demonstrated by the genetic association of the sex phenotype onto the

145   expected chromosome in a narrow window of elevated linkage disequilibrium.

146

147   **Methods**

148

149   **Data availability**

150       All raw sequencing reads and the assembled genome described in this project have

151   been submitted to NCBI under BioProject PRJNA556729. SNP variant sets described below are

152   available through Dryad repository .                                                         **Comment [RE1]:** TBD on acceptance

153

154   **Animal care and sample collection**

155   All animals were reared in compliance with Canadian Council on Animal Care Guidelines, under

156   oversight from the Fisheries and Oceans Canada Pacific Region Animal Care Committee

157   (PRACC). Chum salmon for genome sequencing and assembly and for transcriptome assembly

158   were from Chehalis River Hatchery parents and reared at Fisheries and Oceans Canada in West

159   Vancouver. Chum salmon mitotic gynogen doubled haploids were produced following

160   procedures described by [62]. Briefly, eggs were fertilized with UV-irradiated sperm and

161   pressure shocked (10,000 psi for 5 minutes) in batches at 30 min intervals between 4 and 7

162   hours post-fertilization. One individual from the 7h pressure shock group (Oke142-1, NCBI

163   BioSample: SAMN12367893; Supplementary Table 1) was confirmed to be homozygous for

164   maternal alleles using a panel of 14 microsatellites [49], and was used for genome sequencing

165   and assembly (see below). The individual was euthanized in a bath of 200 mg/L tricaine

166   methanesulfonate (TMS) buffered in 400 mg/L sodium bicarbonate prior to first feeding stage,

167   and stored in ethanol before DNA extraction and whole genome sequencing.

168

169     For transcriptomic data, control Chehalis River Hatchery chum salmon produced from

170     the same parents as Oke142-1 but without UV milt treatment or pressure shock were grown in

171     aerated fresh well water in 200–3700 L tanks and fed hourly as fry and to satiation 3 times daily

172     as parr with stage-appropriate manufactured salmon feed (Skretting Canada Ltd.). At

173     approximately 7 months post-ponding, a single selected chum female (86.9g with a 19.3cm fork

174     length) was euthanized with TMS as above, then rapidly (< three min, PRACC management

175     procedure 3.7) team dissected to harvest 18 tissues (see Supplementary Table 2) for RNA

176     extraction, with an additional tissue (testes) sampled from an juvenile male. All tissues were

177     stored in RNAlater at -20°C until extraction. RNA extractions were performed using the Qiagen

178     RNeasy Mini Kit following the manufacturer's protocol.

179     For individuals used in resequencing, samples were obtained primarily through non-

180     lethal sampling of fin clips or operculum punches from Fisheries and Oceans Canada hatchery

181     brood programs. Additional samples were obtained from archived tissue sets used for genetic

182     stock ID baseline development to supplement the dataset. In total, 59 individuals were utilized

183     in this assessment, with DNA obtained via Qiagen DNeasy Animal tissue kit's following

184     manufacturer's protocol) or phenol/chloroform extractions (following Thermo Fisher Scientific's

185     protocol for genomic DNA preparation [63]. Tissue types, sex, collection dates and locations are

186     summarized in Supplementary Table 3.

187

**188     Genome sequencing and Assembly**

189     DNA was isolated from RNAlater or ethanol preserved tissues using a

190     phenol/chloroform extraction as per Thermo Fisher Scientific's protocol for genomic DNA

191     preparation [63].  Extracted DNA was submitted for genome sequencing across multiple library

192     types, using both Illumina and PacBio sequencing instruments (summarized in Supplementary

193     Table 1: SRA chum Gynogen). Extracted DNA was submitted to the McGill University and

194     Génome Québec Innovation Centre (now the Centre d'expertise et de services Génome

195     Québec) for construction of overlapping (library size estimate = 497 base pairs (bp) and non-

196     overlapping (library size estimate = 620bp) IDT dual-indexed Illumina Shotgun libraries. Each

197     library was sequenced twice on an Illumina HiSeq2500 on RAPID mode PE250. Extracted DNA

198 was also submitted to the McGill University and Génome Québec Innovation Centre for

199 construction of a single library of 10X Chromium linked-reads. Following library construction,

200 the library was sequenced across three lanes of Illumina HiSeqX PE150. Extracted DNA was also

201 submitted to the National Research Council Plant Biotechnology Institute Genome Core for

202 Illumina mate-pair library construction and sequencing. Mate-pair libraries targeting 2-3kb, 4-

203 6kb and 7-12kb were constructed, and sequenced on a lane each of Illumina HiSeq2500 PE125.

204 Finally, extracted DNA was submitted to McGill University and Génome Québec Innovation

205 Centre for construction of a Pacific Biosciences SMRT library using a sheared large insert library

206 type, and the MagBead OneCellPerWell v1 collection protocol. The library was ultimately

207 sequenced across 16 total SMRT cells.

208  Assembly protocols followed successful strategies utilized for Northern Pike e.g., [40, 61,

209 64, 65]. See Supplementary Table 4 for specific parameters to assembly and trimming that were

210 tested. Reads were first trimmed for quality, adapters and minimum length using Trimmomatic

211 [66], and BBmap's FilterByTile was utilized to remove poorly performing portions of the

212 Illumina reads (https://jgi.doe.gov/data-and-tools/bbtools/bb-tools-user-guide/bbmap-guide/;

213 [67]). Allpaths-LG v52488 [68] was utilized with overlapping Illumina overlapping PE250 and

214 Illumina mate-pair libraries using a 3.0 TB memory node on the Compute Canada cluster Cedar.

215 Non-overlapping libraries were also included in two assembly attempts, but ultimately

216 exceeded the memory availability on the node in the MergeNeighbourhoods2 module and

217 were dropped in successful assemblies. Assembly parameters were primarily adjusted for

218 coverage of each of the library types as had been performed in other species; additional

219 modifications were made to read filtering to improve the assemblies.

220  Following Allpaths-LG assembly, scaffolds were passed into PB Jelly 2 v 15.8.24 [69]

221 along with all subreads produced in PacBio sequencing. Nodes on Compute Canada's Cedar

222 cluster were used for all stages, with on-node temp directory and 48 cores used in all steps

223 where allowed. Blasr parameters were `-minMatch 8 -sdpTupleSize 8 -minPctIdentity 75 -bestn

224 1 -nCandidates 10 -nproc 48 -maxScore -500 –noSplitSubreads`. Extraction.py was modified to

225 `MAXGAPHOLD= 1000000` to take advantage of memory available. Collection.py was run with

226 `–m 3`. All other parameters remained default. Finally, the assembly was polished with Pilon

227 [70] using the trimmed paired-end data, aligned to the genome utilizing `bwa mem –M` and

228 default parameters.

229       Scaffolds were ordered and oriented into chromosome representations (i.e.,

230 Pseudomolecules) predominately following the methods described in Christensen et al. (2018)

231 [40]. The sequences underlying the markers for the published chum linkage map from Waples

232 et al. (2016) [34] were aligned to the scaffold assembly utilizing BLAST (-outfmt 6, -word_size

233 48, perc_identity 94, -max_hsps 100, -max_target_seqs 10 -evalue 1E-16). All scaffolds with a

234 link to at least one marker on the map were retained for subsequent pseudomolecule inclusion.

235 Scaffolds were ordered and oriented to the extent allowed by the linkage map, although

236 regions of low recombination limited the effectiveness of the maps alone at this task.

237 Therefore, the sequences underlying the markers for the linkage map were also aligned to a

238 higher contiguity genome of a related species (coho; GCF_002021735.2), and ordering and

239 orientation was further refined based on the conserved synteny between the two species via

240 manual review. Where discrepancies were observed, the chum linkage map was taken as

241 correct to ensure major species-specific rearrangements were captured. Finally,

242 pseudomolecules were aligned to genomes of additional salmonids rainbow trout

243 GCF_002163495.1 [39], Atlantic salmon (GCF_000233375.1) [38], Chinook salmon

244 (GCF_002872995.1) [40] and the non-duplicated outgroup to the salmonids, northern pike

245 (GCF_000721915.3) [71] using Symap v4.2 [72] to ensure linearity was generally conserved, and

246 where it was not, was supported by rearrangements observed in the linkage map.

247       A BUSCO v4.0.2 [73] analysis utilizing the actinopterygii_odb10 dataset and `-m geno –c

248 10 –sp zebrafish` was used to analyze the gene representation within the assembly utilizing the

249 RefSeq maintained assembly: GCF_012931545.1.

250

251 **Gene Annotation**

252       Raw reads for RNA-seq libraries were uploaded into NCBI under BioProject

253 PRJNA556729 for inclusion in the Eukaryotic Genome Annotation pipeline. NEBNext dual-

254 indexed mRNA stranded libraries were constructed from tissues described above by the McGill

255 University and Génome Québec Innovation Centre, and sequenced on a half lane of NovaSeq

256   6000 S4 PE150 (additional libraries in the lane consisted primarily of RNA-seq of Pink and

257   Chinook salmon from related projects). Sequences were uploaded under: SRP216443, with

258   individual accessions: SRR9841162 (Adipose), SRR9841163 (Brain), SRR9841160 (Gill),

259   SRR9841161 (Head Kidney), SRR9841166 (Heart), SRR9841167 (Hindgut), SRR9841164 (Left

260   Eye), SRR9841165 (Liver), SRR9841168 (Lower Jaw), SRR9841169 (Midgut), SRR9841171

261   (Ovary), SRR9841172 (Pituitary), SRR9841170 (Pyloric Caeca), SRR9841174 (Red Muscle Skin),

262   SRR9841176 (Spleen), SRR9841177 (Stomach), SRR9841173 (Testes), SRR9841175 (Upper Jaw

263   Nares), and SRR9841178 (White Muscle).

264

265   **Variant Calling**

266   All individuals sequenced for variant calling (Supplementary Table 3) used Shotgun PCR

267   Free IDT dual-indexed Illumina libraries, produced on a quarter lane of Illumina HiSeqX – library

268   construction and sequencing were performed at the McGill University and Génome Québec

269   Innovation Centre. Raw reads were uploaded to the NCBI BioProject PRJNA556729, with

270   individual accessions listed in the supplementary table.

271   Variant calling followed the best practices pipeline of GATK 3.8 [74–76], and generally

272   followed the methods previously outlined in Christensen et al. (2020) [61]. Raw paired-end

273   reads were aligned to the scaffold-version of the genome (pre-pseudomolecule construction)

274   using bwa (v0.7.17) mem [77] and the `-M` option. Samtools (v1.9) [78] was used to sort and

275   index the alignment files, while Picard (v2.18.9) [79] was utilized with the MarkDuplicates

276   option to identify likely PCR duplicates, and with ReplaceSamHeader to add read group

277   information to the alignment files. GATK's HaplotypeCaller was then used `--genotyping_mode

278   DISCOVERY –emitRefConfidence GVCF` to generate gvcf files, GenotypeGVCFs was used to

279   generate vcf on intervals, and CatVariants was used to concatenate interval files into a single

280   vcf. A training variant set was generated using a hard-filtered subset of the first round of

281   genotyping, utilizing VariantFiltration and the parameters ` --filterExpression "QD < 2.0 || FS >

282   60.0 || MQ < 40.0 || MQRankSum < -12.5 || ReadPosRankSum < -8.0"` as well as the VCFtools

283   (v0.1.14) [80] parameters `--maf0.1 –hwe0.01`. A truth set was generated by overlapping the

284   linkage map SNPs with the hard-filtered training set to obtain SNPs found in both methods.

285  VariantRecalibrator was applied using these sets `-mode SNP -an QD -an MQ -an MQRankSum -

286  an ReadPosRankSum -an FS -an SOR -an InbreedingCoeff`, and ApplyRecalibration run to

287  generate a final SNP set ` --ts_filter_level 99.0`. Finally, as SNP calling began pre-

288  pseudomolecule construction, vcfChromTransfer in the Genomics General repository

289  (https://github.com/simonhmartin/genomics_general; commit: 9d12505) was used to lift over

290  the VCF file based on the NCBI submission AGP file. This lifted-over VCF is included in the

291  accompanying dataset as the "raw SNP" set (referred to as set 1 below).

292      VCFtools v1.14 [80] `--maf 0.05 --max-alleles 2 -- min-alleles 2 --max-missing 0.9 –

293  remove-filtered-all --remove-indels` was used to retain only bi-allelic markers with little missing

294  data and remove the rarest variants (referred to as set 2). The next filter utilized the

295  VCF.Filter.v1.0.py script [61] to remove variants with allelic imbalance `-ab 0.2`, followed by

296  VCFtools to select only the 37 pseudomolecules `--chr` (referred to as set 3). The final filter

297  utilized BCFtools v.1.9 to filter variants for LD in a 20kb window ` +prune, -w 20kb, -l 0.4, -n 2`

298  (referred to as set 4). Finally, VCFtools `–relatedness2` was run to detect closely related

299  individuals. In light of the results, individuals `Oke180104-Fert164` and `Oke171107-D` are

300  recommended to be used cautiously in further analysis as they were deemed most likely to be

301  haploid progeny (expected) and sibling (unexpected) respectively of other individuals in the

302  analysis (can be applied to all sets prior to further analysis using `vcftools –remove`; removed

303  for Figure 3 below, not removed in Figure 2).

304

305  **SNP dataset analyses**

306      SNPhylo [81] was run on the "set 4" dataset, using additional options `-m 0.05 -P

307  Oket_chroms_37_ld0.2 -b –B 1000 -a 37` in order to generate a bootstrapped phylogenetic tree

308  of the chum salmon. Visualization was performed using the Figtree V1.4.4 package

309  (http://tree.bio.ed.ac.uk/software/figtree/). PCA analyses were performed on the same dataset

310  using the R package `SNPrelate`, with full and Canadian-only sample sets plotted– the set 3 is

311  visualized in this work, with all visualization performed using the ggplot2 package [82].

312      The sex phenotypes associated with re-sequencing samples (Supplementary Table 3)

313  were utilized as the basis for a genome-wide association analysis for sex. Utilizing the allele

12

314 balanced SNP set ("set 2" above), VCFtools v1.14 was used to generate input for plink

315 (chromosomes only). An association test was run in PLINK 1.9 [83] using the formatted output

316 data, and resulting Manhattan plot visualized in R [84] using the qqman package [85] . Further

317 visualization of identified SNPs were performed using the Adegenet package [86]. Counts of

318 coverage utilized samtools v1.9 depth, using default parameters to calculate genome wide

319 coverage over each individual *.bam alignment file, and using the `-b` option to restrict the

320 calculation to only the region of the growth hormone 2 gene (*GH2*) demonstrating elevated

321 coverage in the males following a manual review of the alignments using IGV viewer 2.9.4 [87].

322       Duplicated regions, presumably from the Salmon specific 4R duplication event, were

323 identified by alignments using the default settings of SyMap v4.2 [72], using a repeat-masked

324 version of the genome following prior methods [61], by masking WindowMasker-based

325 repetitive regions using ` sed -e '/^>/! s/[[:lower:]]/N/g` from the RefSeq genome. Summary

326 tracks were predominately generated using scripts from [40]: Orientation of the blocks were

327 generated using Analyze_Symap_Block_Orientation.py; percent identity was determined using

328 Analyze_Symap_Linear_Alignments.py; percent identify of repetitive regions identified using

329 Percent_Repeat_Genome_Fasta.py. Linkage map markers from "Map 1" in [34] were aligned to

330 the genome as previously described above using BLAST. Linkage disequilibrium (LD) was

331 interpreted using the `--geno-r2` option in VCFtools [80], and outputting only for those

332 comparisons exceeding `--min-r2 0.5` in order to identify the most highly linked SNPs –

333 summaries were further limited to single chromosomes using the `--chr` option. LD calculations

334 utilized the allele balanced set (set 3) described above. LD track utilized counts of markers in

335 linkage disequilibrium across at least 100kb, and summarized as a log sum per 1 million base

336 pairs. Circos v0.69.9 [88] was utilized to visualize the data tracks described.

337       Heterozygosity analyses followed the same parameters and method as in [61]. Runs of

338 homozygosity were identified from the variants that had been filtered for allele balance using

339 PLINK v1.9 (parameters:—homozyg) [83]. The number of heterozygous genotypes and

340 alternative homozygous genotypes per individual were counted using the same custom script

341 described in the supplementary data of the sockeye genome [61]. Heterozygotes per kbp was

342 calculated as the number of heterozygous genotypes divided by the total nucleotides in the

343   genome (1,853,104,330) multiplied by 1 kbp. The heterozygosity ratio was calculated as the

344   number of heterozygous genotypes divided by the number of alternative homozygous

345   genotypes.

346

347   **Results and Discussion**

348   **Genome Assembly and Annotation**

349         From a raw data set consisting of 59X coverage (110 billion bp) of overlapping 250bp

350   Illumina reads and 60X coverage (114 billion bp) of total mate-pair Illumina reads of three insert

351   sizes (2, 5 and 8kb mean), multiple assembly attempts were performed varying the parameters

352   on read depth as well as read-trimming . Ultimately, three of the attempts resulted in a

353   completed assembly (see Table 1), with the final attempt being the most successful, with a

354   contig N50 of 13.1 kb and a scaffold N50 of 653 kbp. Following AllPaths-LG assembly, contig

355   gaps were filled utilizing PB Suite and 53 billion bp of Pacific Biosciences Sequel long-reads.

356   Following Pilon polishing, utilizing the short insert Illumina libraries, scaffolds were organized

357   into pseudomolecules representing the 37 chromosomes in chum salmon, predominately

358   guided by the publicly available linkage map [34]; ultimately, the linkage map allowed for 70%

359   of the genome assembly to be assigned to a linkage group, slightly lower but approximately

360   equivalent to prior attempts in salmonids using equivalent techniques (e.g., [40, 61]). The final

361   assembly was uploaded to NCBI under BioProject PRJNA556729 and ultimately was included in

362   the RefSeq database as GCF_012931545.1.

363         Busco scores indicate that most of the genome is represented within the family, with

364   results similar to what has been seen in Sockeye, with 85.0% complete (25.1% duplicated), 3.2%

365   fragmented and 11.8% missing. This likely reflects the slightly more fragmented nature of the

366   genome as compared to prior attempts using the same technology in other species. We believe

367   this is most likely due to some minor shearing observed in the DNA utilized for library

368   preparation. We did attempt to use 10X chromium data as part of this assembly process, but,

369   our scaffolding power was negligible – after review, it is likely that DNA shearing noted in the

370   bioanalyzer trace prior to library construction limited the size of the fragments from which to

14

371  generate the linked reads, thus limiting scaffolding power. The raw data from this attempt is

372  included under the BioProject (see Supplementary Table 1), but further attempts would need to

373  use a separate individual in order to increase length of the starting material. Given that

374  sequencing and assembly technology has advanced rapidly since we began this project, it is

375  likely further efforts to improve the genome may benefit from the use of long-read

376  technologies, where incredible advances in contiguity have already been demonstrated in

377  salmonids [44, 89]. Indeed, a long-read assembly for chum salmon is planned by the authors,

378  and will eventually replace this reference, in due course.

379  Following inclusion in the RefSeq database, the genome was annotated utilizing the

380  NCBI Eukaryotic Annotation pipeline, ultimately yielding Annotation Release 100

381  (https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Oncorhynchus_keta/100/) – see

382  Table 2 for a summary. Gene annotation, via chum-specific reads, primarily utilized the 19

383  tissue RNA-seq dataset sequenced as part of this work (see Supplementary Table 2), with

384  additional contribution of sequences from two additional datasets with publicly accessible RNA-

385  seq data [90, 91]. Overall, gene numbers are comparable to other salmonid genomes, and thus

386  likely reflect a relatively complete representation of the coding sequence. It is likely that a

387  future genome that utilized long-reads would result in a slightly increased number of genes (as

388  observed between *Oncorhynchus kisutch* Annotation releases 100 and 101, for example).

389  The resulting assembly is summarized visually in Figure 1. Duplicated regions, as

390  identified via self-alignment using Mummer [92] reflect re-diploidized segments of the genome

391  from the salmon-specific 4R duplication event. There are observations of elevated percent

392  identity on the ends of some chromosomes (Figure 1) that demonstrate partial re-diplodization

393  as in Waples et al. (2016) [34] (e.g., LG05 and LG32), but the effect is not nearly as extensive as

394  that observed in other species. The repetitive elements identified by Window Masker were

395  elevated in regions likely overlapping with centromeres based on synteny with other species for

396  which chromosome arms have been described (Figure 1). Figure 1 also shows Map 1 from

397  Waples et al. (2016) [34] (which can be further visualized in more detail in Supplementary

398  Figure 1 and map 2 in Supplementary Figure 2), and demonstrates the co-linearity of the map

399  with the pseudomolecules. As the maps do contain regions of low-recombination, much of the

15

400   ordering and orientation of the scaffolds into pseudomolecules (but crucially, not the

401   assignment to the pseudomolecule itself) relies heavily in some positions on the long-read and

402   Hi-C based assembly of coho salmon (GCF_002021735.2). Given the extensive conserved

403   synteny and co-linearity between orthologous salmonid chromosome arms demonstrated

404   elsewhere (e.g., [36, 40]), this would appear to be a reasonable approach, and has been part of

405   the development of pseudomolecules for short-read assemblies in salmonids previously (eg.

406   [40]). Regions of the genome with high LD generally overlap with regions of reduced

407   recombination as observed in the linkage map (Figure 1). Further exploration of regions of high

408   LD can be observed in Supplementary Figures 3 and 4.

409        As a final clarification on the assembly presented, we note that pseudomolecules have

410   been named within the publicly available assembly based on the linkage group naming

411   mechanism in Waples et al. (2016), [34] to allow for direct comparison between the two works.

412   However, the authors also note, and are enthusiastic about, the naming convention suggested

413   by Sutherland et. al., [36] to describe chromosomal arms, and indeed the adoption of the

414   system into the grayling genome assembly [42]. We provide here in Table 3 the naming for the

415   pseudomolecules that could be suggested by such a system. While the pattern of fusions do

416   make this system less than ideal, and the resulting chromosome names are somewhat

417   unwieldy, we provide them here as a quick reference and potential guide to re-naming of the

418   linkage groups should such a system continue to prove popular as future assemblies are

419   released. Presenting both names here will hopefully ease future reference, whichever naming

420   scheme ends up being formally adopted in future works.

421

**Population level variation**

423        Given the extensive distribution of chum salmon, attempts were made to maximize

424   geographic distribution of the samples selected within the study. We were able to take

425   advantage of an extensive collection of samples [10] in the archive of the Molecular Genetics

426   Lab (Pacific Biological Station, Fisheries and Oceans Canada), combined with more recent

427   contributions from various Fisheries and Oceans Canada hatchery staff for recent brood. While

428   the collection is focused on British Columbia, the addition of the Japanese samples originating

429    from the Tokushibetsu River on the Island of Hokkaido give a glance at the degree of variation

430    expected across the Pacific. Samples and available metadata are summarized in Supplementary

431    Table 3. In total, 15,372,999 nucleotide variants have been described with this data in the raw

432    dataset, with described filters leaving 8,868,081 in set 2, 2,135,295 in set 3, and 94,080 in set 4.

433    A summary of statistics by individual is given in supplementary table 5 [61]. On average (and

434    ignoring the haploid individual), total lengths of runs-of-homozygosity (ROH) averaged 12.4

435    Mbp [0 - 40.8 Mbp as determined using default parameters]. Heterozygous SNPs per 1kbp

436    averaged 1.47 (Standard Deviation = 0.15), while the heterozygosity ratio averaged 2.23

437    (Standard Deviation = 0.45). Overall, results are relatively similar to what was observed utilizing

438    a parallel analysis in Sockeye salmon, although the overall length of ROH is lower (12.4 Mbp in

439    chum salmon vs. 35.5 Mbp in Sockeye salmon), whereas heterozygous SNPs per 1kbp are

440    increased (1.47 in chum salmon vs. 0.67 in Sockeye salmon), and the heterozygous ratio was

441    approximately equivalent (2.23 in chum salmon vs. 2.21 in sockeye salmon after removing

442    outliers). Deviations below the mean for both heterozygosity calculations were predominately

443    associated with average coverage, implying that depth of sequencing likely impacted to some

444    extent these calculations. Regardless, we demonstrate in chum salmon that there is a general

445    increase in heterozygosity as compared to sockeye salmon, and establishes a comparative

446    metric to be carried through to future comparative analyses in other Pacific salmonids.

447          Analyses of the SNP set resulting from whole genome resequencing (targeted coverage

448    of 15X) should be considered exploratory, as collections were focused on geographic coverage

449    to maximize variants within the catalogues rather than addressing additional questions.

450    Nevertheless, the geographic variation explored allowed us to better understand differentiation

451    among British Columbia locations. To this end, a bootstrapped maximum likelihood tree was

452    constructed using a linkage-disequilibrium thinned SNP-set using SNPhylo [81]. As can be seen

453    in Figure 2, the dendrogram clusters samples by regions similar to past analyses with more

454    comprehensive sampling but using older marker technologies (see above). Samples can be

455    resolved into regions corresponding to descriptions from the comprehensive sampling of

456    Beacham et al. (2009) [10], with individual samples resolvable into Japan – Hokkaido; BC

457    Central Coast (Snootli and Kitimat); BC- Haida Gwaii (Deena Creek); BC – West Coast Vancouver

17

458    Island (Nitinat); BC – Strait of Georgia (Tenderfoot, Big Qualicum, Puntledge); and BC – Lower

459    Fraser (Chilliwack, Inch, Chehalis). However, within clusters from multiple regions, we see a

460    relative lack of resolution to the riverine level. Such observations are supported by Principal

461    Component Analysis (PCA) as well (Figure 3);  however, we do begin to see stronger

462    delineation, possibly from the increased number of variant and dimensions in the PCA analysis.

463    In Figure 3A, we observe differentiation across the Pacific Ocean (best described along PC1),

464    and to a lesser degree geographically across the British Columbia coastline (along PC2). When

465    described regionally, individuals from most populations can easily be resolved when focusing

466    on the British Columbia coastline (Figure 3B), and we are able to see delineation among all

467    collections, except those in the Fraser River Basin. Focusing on the Fraser River Basin sites

468    alone, the pattern is less clustered (Figure 3C), although we do see some separation  from

469    salmon collected in different river systems of the Fraser River drainage.

470        Clustering techniques show that river-level resolution is not always observed. Such

471    results have been noted in the past when considering fishery mixture resolution and describing

472    assignments to region only (for example [9, 93]), but it is worth emphasizing that incomplete

473    resolution among collected populations remains true when considering a relatively

474    comprehensive genome-wide representation of variation. As part of the thinning procedure for

475    SNPhylo, however, by default a significant number of SNPs are removed to increase the speed

476    of the calculation. Alternatively, in the analysis of principal components, with just an LD

477    threshold applied (0.2), a much greater number of SNPs were input into the resulting analysis,

478    and it is likely that the number of SNPs in the end analysis played at least a partial role in the

479    reduced delineation observed in the dendrogram relative to the PCA. While collection level

480    differentiation does emerge in the PCA result, observations on reduced datasets (e.g., by

481    chromosome) greatly inhibited the resolving power of the analysis (supplementary figure 5).

482    Based on the results presented here, is is likely that collection level-specific SNPs could be

483    identified in this dataset that maximize the population differentiation observed genome-wide,

484    and that would further drive differentiation observed in the PCA. However, with such a small

485    sampling size, it is likely that any such discovery would be more a representation of sampling

18

486  depth, and the noise within a set would be high.  However, this dataset is now available, should

487  future researchers need to draw on a pool of potential SNPs from which to develop such assays.

488      Within BC, chum salmon regional groupings are described at the conservation unit (CU)

489  level [94], and it is intriguing to note that there may be substructure to the results observed

490  along those lines in the present analysis. For example, the Tenderfoot hatchery samples in the

491  Howe Sound-Burrard Inlet CU do tend to cluster more strongly, relative to the other collection

492  sites in the adjacent Georgia Strait CU suggesting that a greater sample size may allow recovery

493  of further groupings. However, it is likely that straying, generally described as high in chum

494  salmon, is playing a role in limiting genetic distinctiveness to the level of the CU (or higher)

495  regional groupings. While sampling within the study focused primarily on large hatchery

496  operations, it is also possible we are simply revealing a high degree of variation within each

497  population due to a large effective population size, in which case sufficient additional sampling

498  may coalesce around a mean per population. Still, even within the dataset here, the

499  observation remains that individual population level resolution within a region may begin to be

500  demonstrated with genome-wide representation.

501

**Mapping the sex-determining region**

503      Although limited metadata was collected for individuals sampled beyond geographic

504  locations sampled, we were able to collect phenotypic sex information on hatchery brood

505  samples. Thus, we were able to explore genome wide associations (GWAs) of phenotypic sex.

506  As demonstrated in Figure 4A, two clear peaks were observed with the GWAS: a very strong

507  peak on Linkage Group 15, and another, albeit somewhat weaker, association on Linkage Group

508  3. As shown in Figure 4B, the specific region overlaps with an area of increased linkage

509  disequilibrium on the distal end of LG15 . In Figure 4C, the genotypes for each individual is

510  displayed for the 20 SNPs seen as most associated with sex within the GWAS analysis. LG15 has

511  been previously identified by McKinney et al. (2020), [59] as linked to sex during a RAD-seq

512  based study of chum salmon populations within Alaska. In this prior work, linkage of sex to a

513  particular region of the genome was complicated by two potential factors – a lack of a

514  chromosome-level assembly for chum salmon, and the identification of a putative inversion

515    along the chromosome that resulted in significant patterns of linkage. We utilized the sex-

516    linked RAD loci to position the markers onto the new genome assembly  and observed that

517    while all were indeed placed along Oket_LG15, they appeared to be more dispersed along the

518    chromosome, and were not strongly linked to sex within our geographically distinct sample set

519    (Supplementary Table 7). Within the present study, we observed sex linked to a very narrow

520    region along Oket_LG15; while some noise is observed, the peak is approximately in the 30.8

521    Mbp to 31 Mbp region  and encompasses four annotated genes: potassium/sodium

522    hyperpolarization-activated cyclic nucleotide-gated channel 2-like; E3 ubiquitin-protein ligase

523    RNF126-like; SURP and G-patch domain-containing protein 1-like; and serine/threonine-protein

524    kinase STK11-like. While we do not suggest any of these are the sex-determination gene – as

525    with other Pacific salmonids it is presumed to be sdY [95] – given that the underlying genome

526    assembly is female, this likely represents the approximate region where sdY is inserted on the

527    Y-chromosome, and limited recombination surrounding the region has led to sex-specific

528    markers extending to autosomal-like sequence flanking the insertion. This region (on

529    chromosome 3.2 based on the naming scheme in Sutherland et al., 2016 [36] and Table 3)

530    would appear to be a unique placement thus far in sdY mapping – however, the relatively

531    common observation of sdY on chromosome arm 3.1 (sockeye salmon, coho salmon, lake

532    whitefish; [96] and references therein) does suggest that inter-homeologue transfer between

533    chromosome arms arising from the most recent salmon-specific duplication could be a

534    mechanism for this transfer.

535        The strong secondary peak observed on Linkage group 3 is slightly more confounding

536    and intriguing, as it does not appear to be linked to a known sex-determination orthologue in

537    salmonids [96], and because potential sex-markers appear linked to those on LG-15. While it

538    could be linked to a misplaced contig within the assembly, comparative mapping between

539    additional species did not suggest anything was misplaced based on conserved synteny (data

540    not shown; performed within Symap using default parameters) – if this is the case, it is likely

541    that a future long-read based assembly will correct such a matter. It seems most likely in this

542    case that it represents a repetitive element or otherwise duplicated sequence that is prominent

543    in the Y-specific region but is not present in this female genome; thus, mis-mapping appears to

20

544    occur elsewhere in the genome. A manual review of the region does imply a highly repetitive

545    region, with great differentiation in depths indicative of collapsed repeats. Such mismapping

546    based on collapsed repeats or a lack of sex-specific reference is not uncommon (e.g., as

547    demonstrated in Chinook salmon by mapping of the Y-specific growth hormone pseudogene to

548    the GH2 locus on a different chromosome [97]) and it may be that assembly of a male genome

549    will reveal repeat patterns underlying this unexpected result observed here. There may be

550    additional, more complex reasons based on the observance of multiple sdY regions seen in

551    other species (e.g., Atlantic salmon [98]), although other explanations may be equally likely

552    here. Observations have been made elsewhere that GH-Y, a commonly used proxy for genetic

553    sex in salmonids [99], was found to be missing in males or present in females in some chum

554    salmon populations [100]. While the presented genome is female-based (and thus not

555    predicted to contain GH-Y), observation of relative coverage at the most closely related gene in

556    the genome – GH2 – indeed implied that between 0-5 copies of GH-Y are observed in male

557    individuals, with those males observed to be missing GH-Y being from Kitimat (2x), Snootli (1x)

558    and Tenderfoot (1x): see supplementary table 6. These data do not suggest the phenotypes are

559    mis-identified, however, as inclusion of a Rainbow trout sdY into the alignment phase

560    demonstrated that the presence of sdY matches the phenotype, as would be predicted [95]. No

561    copy-number differences could be interpreted from the sdY alignment unfortunately, as the

562    underlying sequence from trout appeared too differentiated to obtain a reliable estimate of

563    coverage; however, reads were observed aligned to the sequence in all male individuals and

564    not in female individuals in a manual review utilizing IGV viewer. Still, the GH-Y results do

565    indicate that there is variability in the genomic architecture surrounding sdY, and perhaps may

566    indicate that alternate locations within the genome could be influential.  Whatever the

567    underlying genomic architecture of the sex-determination region may be in chum salmon, the

568    result presented here underlines the usefulness and ease of use of the presented SNP dataset

569    and reference genome in mapping a trait of interest to the appropriate chromosome and

570    chromosomal region within the genome.

571

21

**Conclusions**

572
573          The genome assembly for chum salmon represents a relatively complete representation
574     of the chum salmon genome: the first such resource for the species. Contiguity and
575     completeness is likely most affected in regions with high residual tetraploidy or incomplete re-
576     diploidization. While long-read based assemblies (and future sequencing technologies) are
577     likely to generate a more complete picture, the current genome assembly represents a valuable
578     resource for chum salmon on par with those available for Chinook, sockeye, and longstanding
579     assemblies for Atlantic salmon and rainbow trout that allowed a transformation in genomic
580     understanding of these commercially and culturally specific species (e.g., [101]).
581     Complementing the presented genome is a pilot-scale catalogue of variation that provides a
582     genome-wide resource for British Columbian chum salmon populations, and allows for
583     contrasting variation in Western and Eastern Pacific lineages. Such a dataset will be explored
584     further as a resource for SNP genotyping panel expansion, structural variation discovery, or as
585     demonstrated here, in identification of the chromosome and position most likely to contain the
586     sex-determination gene in chum salmon.
587

22

**References**

1. Helfield JM, Naiman RJ. Keystone Interactions: Salmon and Bear in Riparian Forests of Alaska. Ecosystems. 2006;9:167–80.

2. Salo EO. Life History of Chum Salmon (Oncorhynchus keta). In: Groot C, Margolis L, editors. Pacific salmon life histories. Vancouver: UBC Press; 1991.

3. Bakkala RG. Synopsis of Biological Data on the Chum Salmon, Oncorhynchus Keta (Walbaum) 1792. U.S. Fish and Wildlife Service; 1970.

4. Fredin RA, Major RL, Bakkala RG, Tanonaka GK. Pacific salmon and the high seas salmon fisheries of Japan. 1977.

5. Behnke R. Trout and salmon of north america. Free Press; 2010.

6. Gislason G, Lam E, Knapp G, Guettabi M. Economic Impacts of Pacific Salmon Fisheries. Pacific Salmon Commission. University of Alaska Anchorage Institute of Social & Economic Research.

7. Keefer ML, Caudill CC. Homing and straying by anadromous salmonids: a review of mechanisms and rates. Reviews in Fish Biology and Fisheries. 2014;24:333–68.

8. Quinn TP, Stewart IJ, Boatright CP. Experimental evidence of homing to site of incubation by mature sockeye salmon, Oncorhynchus nerka. Animal Behaviour. 2006;72:941–9.

9. Beacham T, Sato S, Urawa S, Le K, Wetklo M. Population structure and stock identification of chum salmon *Oncorhynchus keta* from Japan determined by microsatellite DNA variation. Fisheries Science. 2008;74:983–94.

10. Beacham TD, Candy JR, Le KD, Wetklo M. Population structure of chum salmon (Oncorhynchus keta) across the Pacific Rim, determined from microsatellite analysis. Fishery Bulletin. 2009;107:244–60.

11. Olsen JB, Flannery BG, Beacham TD, Bromaghin JF, Crane PA, Lean CF, et al. The influence of hydrographic structure and seasonal run timing on genetic diversity and isolation-by-distance in chum salmon (Oncorhynchus keta). Can J Fish Aquat Sci. 2008;65:2026–42.

12. Small MP, Frye AE, Von Bargen JF, Young SF. Genetic Structure of Chum Salmon (Oncorhynchus keta) Populations in the Lower Columbia River: Are Chum Salmon in Cascade Tributaries Remnant Populations? Conservation Genetics. 2006;7:65–78.

13. Small MP, Rogers Olive SD, Seeb LW, Seeb JE, Pascal CE, Warheit KI, et al. Chum Salmon Genetic Diversity in the Northeastern Pacific Ocean Assessed with Single Nucleotide Polymorphisms (SNPs): Applications to Fishery Management. North American Journal of Fisheries Management. 2015;35:974–87.

631  14. Crow KD. What Is the Role of Genome Duplication in the Evolution of Complexity and
632  Diversity? Molecular Biology and Evolution. 2006;23:887–92.

633  15. Ohno S. Evolution by Gene Duplication. Berlin, Heidelberg: Springer Berlin Heidelberg; 1970.

634  16. Otto SP, Whitton J. Polyploid Incidence and Evolution. Annu Rev Genet. 2000;34:401–37.

635  17. Taylor JS, Raes J. Duplication and Divergence: The Evolution of New Genes and Old Ideas.
636  Annu Rev Genet. 2004;38:615–43.

637  18. Sankoff D, Zheng C. Whole Genome Duplication in Plants: Implications for Evolutionary
638  Analysis. In: Setubal JC, Stoye J, Stadler PF, editors. Comparative Genomics. New York, NY:
639  Springer New York; 2018. p. 291–315.

640  19. Wolfe KH, Shields DC. Molecular evidence for an ancient duplication of the entire yeast
641  genome. Nature. 1997;387:708–13.

642  20. Kellis M, Birren BW, Lander ES. Proof and evolutionary analysis of ancient genome
643  duplication in the yeast Saccharomyces cerevisiae. Nature. 2004;428:617–24.

644  21. Kenny NJ, Chan KW, Nong W, Qu Z, Maeso I, Yip HY, et al. Ancestral whole-genome
645  duplication in the marine chelicerate horseshoe crabs. Heredity. 2016;116:190–9.

646  22. Schwager EE, Sharma PP, Clarke T, Leite DJ, Wierschin T, Pechmann M, et al. The house
647  spider genome reveals an ancient whole-genome duplication during arachnid evolution. BMC
648  Biol. 2017;15:62.

649  23. Putnam NH, Butts T, Ferrier DEK, Furlong RF, Hellsten U, Kawashima T, et al. The amphioxus
650  genome and the evolution of the chordate karyotype. Nature. 2008;453:1064–71.

651  24. Dehal P, Boore JL. Two Rounds of Whole Genome Duplication in the Ancestral Vertebrate.
652  PLoS Biol. 2005;3:e314.

653  25. Taylor JS, Van de Peer Y, Braasch I, Meyer A. Comparative genomics provides evidence for
654  an ancient genome duplication event in fish. Phil Trans R Soc Lond B. 2001;356:1661–79.

655  26. Taylor JS. Genome Duplication, a Trait Shared by 22,000 Species of Ray-Finned Fish.
656  Genome Research. 2003;13:382–90.

657  27. Hoegg S, Brinkmann H, Taylor JS, Meyer A. Phylogenetic Timing of the Fish-Specific Genome
658  Duplication Correlates with the Diversification of Teleost Fish. J Mol Evol. 2004;59:190–203.

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

661  29. Macqueen DJ, Johnston IA. A well-constrained estimate for the timing of the salmonid
662  whole genome duplication reveals major decoupling from species diversification. Proc R Soc B.
663  2014;281:20132881.

664  30. Van de Peer Y, Mizrachi E, Marchal K. The evolutionary significance of polyploidy. Nature
665  Reviews Genetics. 2017;18:411–24.

666  31. Jones SR, Fast MD, Johnson SC, Groman DB. Differential rejection of salmon lice by pink and
667  chum salmon: disease consequences and expression of proinflammatory genes. Dis Aquat
668  Organ. 2007;75:229–38.

669  32. Sutherland BJ, Koczka KW, Yasuike M, Jantzen SG, Yazawa R, Koop BF, et al. Comparative
670  transcriptomics of Atlantic Salmo salar, chum Oncorhynchus keta and pink salmon O. gorbuscha
671  during infections with salmon lice Lepeophtheirus salmonis. BMC Genomics. 2014;15:200.

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

675  34. Waples RK, Seeb LW, Seeb JE. Linkage mapping with paralogs exposes regions of residual
676  tetrasomic inheritance in chum salmon ( *Oncorhynchus keta* ). Mol Ecol Resour. 2016;16:17–28.

677  35. Robertson FM, Gundappa MK, Grammes F, Hvidsten TR, Redmond AK, Lien S, et al. Lineage-
678  specific rediploidization is a mechanism to explain time-lags between genome duplication and
679  evolutionary diversification. Genome Biology. 2017;18:111.

680  36. Sutherland BJG, Gosselin T, Normandeau E, Lamothe M, Isabel N, Audet C, et al. Salmonid
681  Chromosome Evolution as Revealed by a Novel Method for Comparing RADseq Linkage Maps.
682  Genome Biology and Evolution. 2016;8:3600–17.

683  37. Berthelot C, Brunet F, Chalopin D, Juanchich A, Bernard M, Noël B, et al. The rainbow trout
684  genome provides novel insights into evolution after whole-genome duplication in vertebrates.
685  Nat Commun. 2014;5:3657.

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

688  39. Pearse DE, Barson NJ, Nome T, Gao G, Campbell MA, Abadía-Cardoso A, et al. Sex-
689  dependent dominance maintains migration supergene in rainbow trout. Nat Ecol Evol.
690  2019;3:1731–42.

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

693   41. Narum SR, Di Genova A, Micheletti SJ, Maass A. Genomic variation underlying complex life-
694   history traits revealed by genome sequencing in Chinook salmon. Proc R Soc B.
695   2018;285:20180935.

696   42. Sävilammi T, Primmer CR, Varadharajan S, Guyomard R, Guiguen Y, Sandve SR, et al. The
697   Chromosome-Level Genome Assembly of European Grayling Reveals Aspects of a Unique
698   Genome Evolution Process Within Salmonids. G3. 2019;9:1283–94.

699   43. Varadharajan S, Sandve SR, Gillard GB, Tørresen OK, Mulugeta TD, Hvidsten TR, et al. The
700   Grayling Genome Reveals Selection on Gene Expression Regulation after Whole-Genome
701   Duplication. Genome Biology and Evolution. 2018;10:2785–800.

702   44. De-Kayne R, Zoller S, Feulner PGD. A *de novo* chromosome-level genome assembly of
703   *Coregonus sp. "Balchen"* : one representative of the Swiss Alpine whitefish radiation. preprint.
704   Genomics; 2019.

705   45. Beacham TD, Wallace CG, Jonsen K, Sutherland BJG, Gummer C, Rondeau EB. Estimation of
706   Conservation Unit and population contribution to Chinook salmon mixed-stock fisheries in
707   British Columbia, Canada using direct DNA sequencing for single nucleotide polymorphisms.
708   Can J Fish Aquat Sci. 2021. https://doi.org/10.1139/cjfas-2020-0462.

709   46. Beacham TD, Wallace C, Jonsen K, McIntosh B, Candy JR, Rondeau EB, et al. Accurate
710   estimation of conservation unit contribution to coho salmon mixed-stock fisheries in British
711   Columbia, Canada, using direct DNA sequencing for single nucleotide polymorphisms. Can J Fish
712   Aquat Sci. 2020;77:1302–15.

713   47. Phelps SR, LeClair LL, Young S, Blankenship HL. Genetic Diversity Patterns of Chum Salmon in
714   the Pacific Northwest. Can J Fish Aquat Sci. 1994;51:65–83.

715   48. Seeb LW, Crane PA. High Genetic Heterogeneity in Chum Salmon in Western Alaska, the
716   Contact Zone between Northern and Southern Lineages. Transactions of the American Fisheries
717   Society. 1999;128:58–87.

718   49. Beacham TD, Spilsted B, Le KD, Wetklo M. Population structure and stock identification of
719   chum salmon (Oncorhynchus keta) from British Columbia determined with microsatellite DNA
720   variation. Can J Zool. 2008;86:1002–14.

721   50. Smith CT, Seeb LW. Number of Alleles as a Predictor of the Relative Assignment Accuracy of
722   Short Tandem Repeat (STR) and Single-Nucleotide-Polymorphism (SNP) Baselines for Chum
723   Salmon. Transactions of the American Fisheries Society. 2008;137:751–62.

724   51. Smith CT, Baker J, Park L, Seeb LW, Elfstrom C, Abe S, et al. Characterization of 13 single
725   nucleotide polymorphism markers for chum salmon: PRIMER NOTE. Molecular Ecology Notes.
726   2005;5:259–62.

727    52. Smith CT, Elfstrom CM, Seeb LW, Seeb JE. Use of sequence data from rainbow trout and
728    Atlantic salmon for SNP detection in Pacific salmon: SNPs IN PACIFIC SALMON. Molecular
729    Ecology. 2005;14:4193–203.

730    53. Elfstrom CM, Smith CT, Seeb LW. Thirty-eight single nucleotide polymorphism markers for
731    high-throughput genotyping of chum salmon. Mol Ecol Notes. 2007;7:1211–5.

732    54. Seeb JE, Pascal CE, Grau ED, Seeb LW, Templin WD, Harkins T, et al. Transcriptome
733    sequencing and high-resolution melt analysis advance single nucleotide polymorphism
734    discovery in duplicated salmonids: PERMANENT GENETIC RESOURCES ARTICLE. Molecular
735    Ecology Resources. 2011;11:335–48.

736    55. Petrou EL, Hauser L, Waples RS, Seeb JE, Templin WD, Gomez-Uchida D, et al. Secondary
737    contact and changes in coastal habitat availability influence the nonequilibrium population
738    structure of a salmonid ( *Oncorhynchus keta* ). Mol Ecol. 2013;22:5848–60.

739    56. Small M, Warheit K, Pascal C, Seeb L, Ruff C, Zischke J, et al. Chum Salmon Southern Area
740    Genetic Baseline Enhancement Part 1 and Part 2: Amplicon Development, Expanded Baseline
741    Collections, and Genotyping.

742    57. Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, et al. Rapid SNP Discovery
743    and Genetic Mapping Using Sequenced RAD Markers. PLoS ONE. 2008;3:e3376.

744    58. Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA. Rapid and cost-effective
745    polymorphism identification and genotyping using restriction site associated DNA (RAD)
746    markers. Genome Research. 2007;17:240–8.

747    59. McKinney G, McPhee MV, Pascal C, Seeb JE, Seeb LW. Network Analysis of Linkage
748    Disequilibrium Reveals Genome Architecture in Chum Salmon. G3: Genes|Genomes|Genetics.
749    2020;10:1553.

750    60. Gao G, Nome T, Pearse DE, Moen T, Naish KA, Thorgaard GH, et al. A New Single Nucleotide
751    Polymorphism Database for Rainbow Trout Generated Through Whole Genome Resequencing.
752    Front Genet. 2018;9:147.

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

756    62. Quillet E, Garcia P, Guyomard R. Analysis of the production of all homozygous lines of
757    rainbow trout by gynogenesis. J Exp Zool. 1991;257:367–74.

758    63. Genomic DNA Preparation from RNAlaterTM Preserved Tissues—CA [Internet].
759    https://www.thermofisher.com/ca/en/home/references/protocols/nucleic-acid-purification-
760    and-analysis/rna-protocol/genomic-dna-preparation-from-rnalater-preserved-tissues.html.
761    Accessed 18 Feb 2021.

762 64. Christensen KA, Rondeau EB, Minkley DR, Leong JS, Nugent CM, Danzmann RG, et al. The
763 Arctic charr (Salvelinus alpinus) genome and transcriptome assembly. PLoS ONE.
764 2018;13:e0204076.

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

768 66. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data.
769 Bioinformatics. 2014;30:2114–20.

770 67. Marić J. Long Read RNA-seq Mapper. Master Thesis. University of Zagreb; 2015.

771 68. Gnerre S, MacCallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, et al. High-quality
772 draft assemblies of mammalian genomes from massively parallel sequence data. Proceedings of
773 the National Academy of Sciences. 2011;108:1513–8.

774 69. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the Gap: Upgrading
775 Genomes with Pacific Biosciences RS Long-Read Sequencing Technology. PLoS ONE.
776 2012;7:e47768.

777 70. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An Integrated
778 Tool for Comprehensive Microbial Variant Detection and Genome Assembly Improvement.
779 PLOS ONE. 2014;9:e112963.

780 71. Johnson HA, Rondeau EB, Minkley DR, Leong JS, Whitehead J, Despins CA, et al. Population
781 genomics of North American northern pike: variation and sex-specific signals from a
782 chromosome-level, long read genome assembly. bioRxiv. 2020;:2020.06.18.157701.

783 72. Soderlund C, Bomhoff M, Nelson WM. SyMAP v3.4: a turnkey synteny system with
784 application to plant genomes. Nucleic Acids Research. 2011;39:e68–e68.

785 73. Seppey M, Manni M, Zdobnov EM. BUSCO: Assessing Genome Assembly and Annotation
786 Completeness. In: Kollmar M, editor. Gene Prediction: Methods and Protocols. New York, NY:
787 Springer New York; 2019. p. 227–45.

788 74. Poplin R, Ruano-Rubio V, DePristo MA, Fennell TJ, Carneiro MO, Van der Auwera GA, et al.
789 Scaling accurate genetic variant discovery to tens of thousands of samples. bioRxiv.
790 2018;:201178.

791 75. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A framework for
792 variation discovery and genotyping using next-generation DNA sequencing data. Nat Genet.
793 2011;43:491–8.

794  76. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, del Angel G, Levy-Moonshine A, et al.
795  From FastQ Data to High-Confidence Variant Calls: The Genome Analysis Toolkit Best Practices
796  Pipeline. Current Protocols in Bioinformatics. 2013;43:11.10.1-11.10.33.

797  77. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. 2013.

798  78. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
799  Alignment/Map format and SAMtools. Bioinformatics. 2009;25:2078–9.

800  79. Picard toolkit. Broad Institute; 2019.

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

803  81. Lee T-H, Guo H, Wang X, Kim C, Paterson AH. SNPhylo: a pipeline to construct a
804  phylogenetic tree from huge SNP data. BMC Genomics. 2014;15:162.

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

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

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

810  85. Turner SD. qqman: an R package for visualizing GWAS results using Q-Q and manhattan
811  plots. Journal of Open Source Software. 2018;3:731.

812  86. Jombart T. adegenet: a R package for the multivariate analysis of genetic markers.
813  Bioinformatics. 2008;24:1403–5.

814  87. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, et al.
815  Integrative genomics viewer. Nature Biotechnology. 2011;29:24–6.

816  88. Krzywinski M, Schein J, Birol İ, Connors J, Gascoyne R, Horsman D, et al. Circos: An
817  information aesthetic for comparative genomics. Genome Research. 2009;19:1639–45.

818  89. Gao G, Magadan S, Waldbieser GC, Youngblood RC, Wheeler PA, Scheffler BE, et al. A long
819  reads-based de-novo assembly of the genome of the Arlee homozygous line reveals
820  chromosomal rearrangements in rainbow trout. G3 Genes|Genomes|Genetics. 2021.
821  https://doi.org/10.1093/g3journal/jkab052.

822  90. Palstra AP, Fukaya K, Chiba H, Dirks RP, Planas JV, Ueda H. The Olfactory Transcriptome and
823  Progression of Sexual Maturation in Homing Chum Salmon Oncorhynchus keta. PLOS ONE.
824  2015;10:e0137404.

825    91. Tatara Y, Kakizaki I, Kuroda Y, Suto S, Ishioka H, Endo M. Epiphycan from salmon nasal
826    cartilage is a novel type of large leucine-rich proteoglycan. Glycobiology. 2013;23:993–1003.

827    92. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and
828    open software for comparing large genomes. Genome Biology. 2004;5:R12.

829    93. SEEB LW, TEMPLIN WD, SATO S, ABE S, WARHEIT K, PARK JY, et al. Single nucleotide
830    polymorphisms across a species' range: implications for conservation studies of Pacific salmon.
831    Molecular Ecology Resources. 2011;11:195–217.

832    94. Fishery & Assessment Data Section, Pacific Biological Station. Chum Salmon (Oncorhynchus
833    keta) Conservation Units, Sites & Status. 2017.

834    95. Yano A, Nicol B, Jouanno E, Quillet E, Fostier A, Guyomard R, et al. The sexually dimorphic
835    on the Y-chromosome gene (sdY) is a conserved male-specific Y-chromosome sequence in many
836    salmonids. Evol Appl. 2013;6:486–96.

837    96. Sutherland BJG, Rico C, Audet C, Bernatchez L. Sex Chromosome Evolution, Heterochiasmy,
838    and Physiological QTL in the Salmonid Brook Charr Salvelinus fontinalis. G3 (Bethesda).
839    2017;7:2749–62.

840    97. Micheletti SJ, Narum SR. Utility of pooled sequencing for association mapping in nonmodel
841    organisms. Molecular Ecology Resources. 2018;18:825–37.

842    98. Eisbrenner WD, Botwright N, Cook M, Davidson EA, Dominik S, Elliott NG, et al. Evidence for
843    multiple sex-determining loci in Tasmanian Atlantic salmon (Salmo salar). Heredity.
844    2014;113:86–92.

845    99. Devlin RH, Biagi CA, Smailus DE. Genetic mapping of Y-chromosomal DNA markers in Pacific
846    salmon. Genetica. 2001;111:43–58.

847    100. Muttray AF, Sakhrani D, Smith JL, Nakayama I, Davidson WS, Park L, et al. Deletion and
848    Copy Number Variation of Y-Chromosomal Regions in Coho Salmon, Chum Salmon, and Pink
849    Salmon Populations. Transactions of the American Fisheries Society. 2017;146:240–51.

850    101. Bobe J, Marandel L, Panserat S, Boudinot P, Berthelot C, Quillet E, et al. 2 - The rainbow
851    trout genome, an important landmark for aquaculture and genome evolution. In: MacKenzie S,
852    Jentoft S, editors. Genomics in Aquaculture. San Diego: Academic Press; 2016. p. 21–43.

853

854    **Tables**

855

856    Table 1: Assembly results for Allpaths and PBSuite based assemblies performed.

857

858     Table 2: Summary of Annotation Release 100 from the NCBI Eukaryotic Annotation pipeline.

859     See https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Oncorhynchus_keta/100/ for

860     more details.

861

862     Table 3: Pike-like chromosome naming for the chum salmon pseudomolecules described in this

863     work, based on Sutherland et al. (2016) [36]

864

865     **Figures**

866

867     Figure 1: Circos plot of the chum salmon genome GCF_012931545.1. Inner ribbons demonstrate

868     ohnologous regions (regions duplicated at the salmon-specific genome duplication event).

869     Working in to out, Track A describes the average percent identity between the duplicated

870     regions, in 1 Mbp bins. Track B describes the average percent identity in the chromosomes, in 1

871     Mbp bins. Track C describes the relationship to "Map 1" chum linkage map from Waples et al.

872     (2016) [34]. Track D describes SNPs demonstrating elevated LD (R-squared >= 0.5) and >= 100

873     kb apart, demonstrated as a log10 based count, in 1Mbp bins.

874

875     Figure 2: Dendrogram produced by SNPhylo, utilizing set 4 SNP data described in the text.

876     Values at nodes indicate bootstrapping. Samples are coloured by geographic region.

877

878     Figure 3: Principal component analyses performed on set 3 SNPs described in the text, using

879    SNPrelate and plotted in ggplot2. Samples are coloured by collection and displayed in the

880    legend. A) the full dataset (all samples) are presented. B), Japanese samples are removed from

881    the analysis. C), the collections are reduced solely to the collections within the Fraser River

882    drainage.

883

884    Figure 4: Association of the phenotypic sex to the genome utilizing SNP variant set 1. A) the

885    results of the GWAS are presented, with Bonferroni-adjusted p-values shown at the 5% level

886    (blue line) and 1% (orange line) levels. B) The SNPs with R-squared greater than 0.5 are

887    counted, and plotted to show relationship of distance between SNPs being measured, for the

888    region flanking the signal on Oket_LG15. C) The genotypes for each individual is displayed for

889    the 20 SNPs seen as most associated within the GWAS analysis, with homozygous reference in

890    blue, heterozygous in purple, homozygous alternate in red, and missing genotypes in white.

891    Samples are sorted to group males, females and unknowns (Japanese samples—most likely

892    females).

893

894    **Supplementary Data**

895

896    Supplementary Table 1: Biosample and SRA data for individual chum used in generating the

897    genome assembly.

898

899    Supplementary Table 2: Biosample and SRA data for individual chum used in generating the

900    Illumina RNA-seq data.

901

32

902 Supplementary Table 3: Biosample and SRA data for individual chum used in generating the Re-

903 sequencing data.

904

905 Supplementary Table 4: Allpaths-LG parameters explored in attempting to obtain the highest

906 contiguity assemblies.

907

908 Supplementary Table 5: Heterozygosity metrics by individual are described. Includes counts of

909 missing genotypes, Homozygous Reference and Alternate genotypes, Heterozygous genotypes,

910 average depth at called sites, the mean count of heterozygous SNPs per kbp, the ratio of

911 Heterozygous genotypes to Homozygous alternate, and the total length of runs-of

912 homozygosity as determined from PLINK using default parameters.

913

914 Supplementary Table 6: Depth of coverage across the alignments, and at the GH2 locus to

915 approximate the count of GH-Y copies in each individual. GH2 is used for this calculation due to

916 the lack of GH-Y in the reference genome, and therefore the alignment of GH-Y to the closest

917 homologue.

918

919 Supplementary Table 7: Placement of SNPs associated with phenotypic sex from McKinney et

920 al. [59] in Alaskan chum populations onto the current reference genome.

921

922 Supplementary Figure 1: Plotting the association between Linkage groups in Waples et al.

923 (2016), [34] map 1, and the reference genome assembly presented in this work.

924

925 Supplementary Figure 2: Plotting the association between Linkage groups in Waples et al.

926 (2016), [34] map 2, and the reference genome assembly presented in this work.

927

928 Supplementary Figure 3: Plotting the linkage disequilibrium along each chromosome. SNPs are

929 only displayed if R-squared is greater than 0.5, and is plotted as a count of SNPs.

930

931    Supplementary Figure 4: Plotting the linkage disequilibrium along each chromosome. SNPs are

932    only displayed if R-squared is greater than 0.5, and each SNP is plotted by R-squared value.
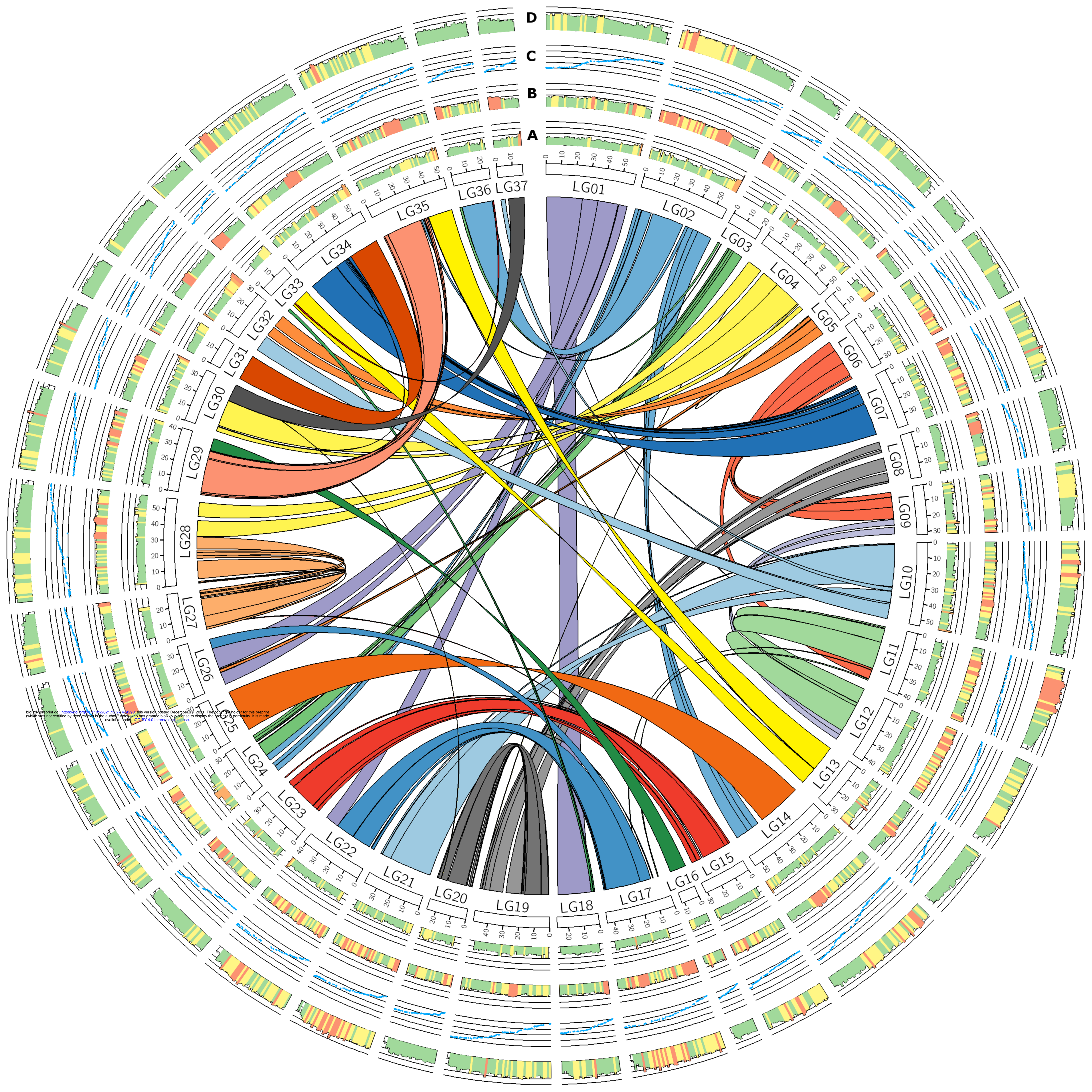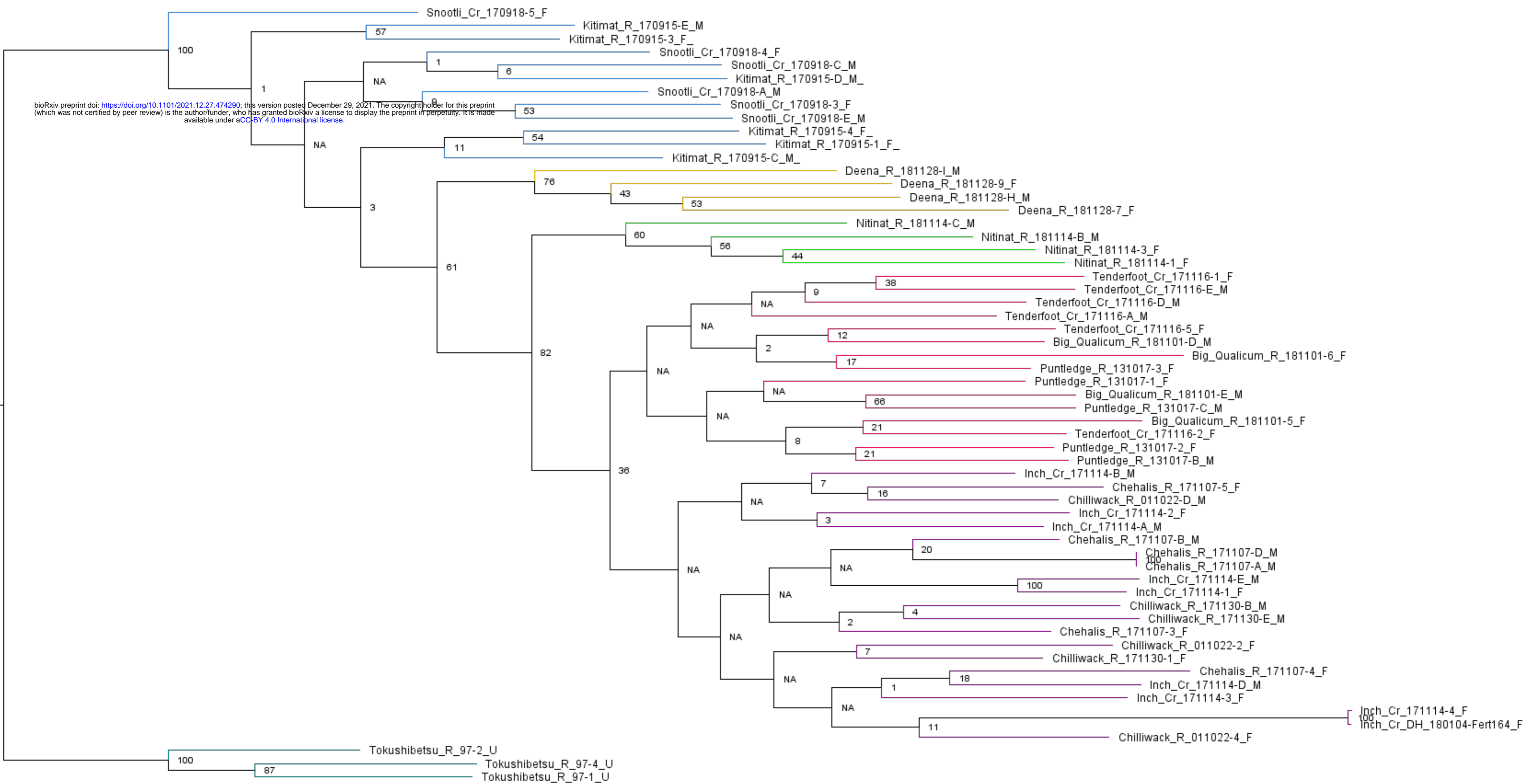
933

934    Supplementary Figure 5: Principal component analyses performed on set 3 SNPs described in

935    the text, using SNPrelate and plotted in ggplot2 and reduced to only query LG15. Samples are

936    coloured by collection, displayed in the legend. In panel A, the full dataset (all samples) are

937    presented. In panel B, Japanese samples are removed from the analysis. In panel C, the samples

938    are coloured by collection site rather than by region.

|  | Assembly size (Contigs) | Assembly size (Scaffolds) | Scaffold N50 | Contig N50 |
|---|---|---|---|---|
| **Allpaths** | 1471097779 | 1813373414 | 653 | 13.1 |
| **PBSuite** | 1,766,907,823 | 1,852,809,593 | 665,581 | 52,191 |

| Feature | Annotation Release 100 |
|---|---|
| Genes and pseudogenes | 45643 |
| protein-coding | 36325 |
| non-coding | 6205 |
| transcribed pseudogenes | 222 |
| non-transcribed pseudogenes | 2821 |
| genes with variants | 14102 |
| immunoglobulin/T-cell receptor gene segments | 70 |

| Linkage Group | Accession | Alternate Naming | |
|---|---|---|---|
| Oket_LG01 | NC_050106 | 1.2-8.2 | |
| Oket_LG02 | NC_050107 | 6.1-2.2 | |
| Oket_LG03 | NC_050108 | | 5.1 |
| Oket_LG04 | NC_050109 | 21.1-4.2 | |
| Oket_LG05 | NC_050110 | | 11.2 |
| Oket_LG06 | NC_050111 | | 14.2 |
| Oket_LG07 | NC_050112 | | 16.1 |
| Oket_LG08 | NC_050113 | | 18.2 |
| Oket_LG09 | NC_050114 | | 14.1 |
| Oket_LG10 | NC_050115 | 17.1-9.2 | |
| Oket_LG11 | NC_050116 | | 7.2 |
| Oket_LG12 | NC_050117 | | 7.1 |
| Oket_LG13 | NC_050118 | 24.1-23.1 | |
| Oket_LG14 | NC_050119 | 13.1-2.1 | |
| Oket_LG15 | NC_050120 | | 3.2 |
| Oket_LG16 | NC_050121 | | 25.1 |
| Oket_LG17 | NC_050122 | | 19.1 |
| Oket_LG18 | NC_050123 | | 1.1 |
| Oket_LG19 | NC_050124 | 10.2-18.1 | |
| Oket_LG20 | NC_050125 | | 10.1 |
| Oket_LG21 | NC_050126 | | 17.2 |
| Oket_LG22 | NC_050127 | | 19.2 |
| Oket_LG23 | NC_050128 | | 3.1 |
| Oket_LG24 | NC_050129 | | 5.2 |
| Oket_LG25 | NC_050130 | | 13.2 |
| Oket_LG26 | NC_050131 | | 8.1 |
| Oket_LG27 | NC_050132 | | 12.1 |
| Oket_LG28 | NC_050133 | 12.2-21.2 | |
| Oket_LG29 | NC_050134 | 15.2-25.2 | |
| Oket_LG30 | NC_050135 | 4.1-22.2 | |
| Oket_LG31 | NC_050136 | | 20.1 |
| Oket_LG32 | NC_050137 | 9.1-11.1 | |
| Oket_LG33 | NC_050138 | | 23.2 |
| Oket_LG34 | NC_050139 | 16.2-20.2 | |
| Oket_LG35 | NC_050140 | 15.1-24.2 | |
| Oket_LG36 | NC_050141 | | 6.2 |
| Oket_LG37 | NC_050142 | | 22.1 |

A)

B) Oket_LG15

C)