# The structural variation landscape in 492 Atlantic salmon genomes

Alicia C. Bertolotti [1,2], Ryan M. Layer [3,4], Manu Kumar Gundappa [2], Michael D. Gallagher [2], Ege Pehlivanoglu [2], Torfinn Nome [5], Diego Robledo [2], Matthew P. Kent [5], Line L. Røsæg [5], Matilde M. Holen [5], Teshome D. Mulugeta [5], Thomas J. Ashton [6], Kjetil Hindar [7], Harald Sægrov [8], Bjørn Florø-Larsen [9], Jaakko Erkinaro [10], Craig R. Primmer [11], Louis Bernatchez [12], Samuel A.M. Martin [1], Ian A. Johnston [6], Simen R. Sandve [5], Sigbjørn Lien [5] *, Daniel J. Macqueen [2] *

* Corresponding authors:

Daniel J. Macqueen (daniel.macqueen@roslin.ed.ac.uk)

Sigbjørn Lien (sigbjorn.lien@nmbu.no)

Affiliations:

[1] *School of Biological Sciences, University of Aberdeen, Tillydrone Avenue, Aberdeen, UK*

[2] *The Roslin Institute and Royal (Dick) School of Veterinary Studies, University of Edinburgh, Edinburgh, UK*

[3] *BioFrontiers Institute, University of Colorado, Boulder, CO, USA*

[4] *Department of Computer Science, University of Colorado, Boulder, CO, USA*

[5] *Centre for Integrative Genetics, Department of Animal and Aquacultural Sciences, Faculty of Biosciences, Norwegian University of Life Sciences, Ås, Norway*

[6] *Xelect Ltd, Horizon House, St Andrews, Scotland, United Kingdom*

[7] *Norwegian Institute for Nature Research (NINA), P.O. Box 5685 Torgarden, NO, 7485, Trondheim, Norway*

[8] *Rådgivende Biologer AS, Bergen, Norway*

[9] *Norwegian Veterinary Institute, P.O. Box 750 Sentrum, 0106, Oslo, Norway*

[10] *Natural Resources Institute Finland (Luke), Oulu, P.O.Box 413 FI-90014, Finland*

[11] *Institute for Biotechnology, University of Helsinki, Helsinki, Finland*

[12] *Institut de Biologie Intégrative et des Systèmes (IBIS) Pavillon Charles-Eugène Marchand, Université Laval Québec QC Canada*

1 **Abstract**

2 Structural variants (SVs) are a major source of genetic and phenotypic variation, but remain challenging to

3 accurately type and are hence poorly characterized in most species. We present an approach for reliable SV

4 discovery in non-model species using whole genome sequencing and report 15,483 high-confidence SVs in

5 492 Atlantic salmon (*Salmo salar* L.) sampled from a broad phylogeographic distribution. These SVs

6 recover population genetic structure with high resolution, include an active DNA transposon, widely affect

7 functional features, and overlap more duplicated genes retained from an ancestral salmonid

8 autotetraploidization event than expected. Changes in SV allele frequency between wild and farmed fish

9 indicate polygenic selection on behavioural traits during domestication, targeting brain-expressed synaptic

10 networks linked to neurological disorders in humans. This study offers novel insights into the role of SVs

11 in genome evolution and the genetic architecture of domestication traits, along with resources supporting

12 reliable SV discovery in non-model species.

3

## Main

13

14 Modern genetics remains primarily focused on single nucleotide polymorphism (SNP) analyses, with a

15 growing recognition of the importance of larger structural variants (SVs) including inversions, insertions,

16 deletions and copy number variations (CNVs) (defined here as variants ≥100 bp), among others[1]. SVs

17 affect a larger proportion of bases in human genomes than SNPs[4], are not always reliably tagged by SNPs[5],

18 more frequently have regulatory impacts[6], and have been shown to alter the structure, presence, number,

19 dosage, and regulation of many genes[1]. Nonetheless, SVs remain challenging to accurately type using

20 whole genome sequence data[2-3], limiting our understanding of their biological roles and exploitation as

21 genetic markers. Consequently, there is a need for reliable SV detection approaches to fully exploit the fast-

22 accumulating genome sequencing datasets in both model and non-model species, allowing for more

23 complete genetics investigations. Many tools exist for SV discovery using short-read sequencing data, but

24 all suffer from high false discovery rates (10-89%)[2,3,7]. This poses a challenge for truly *de novo* SV

25 detection in previously unstudied species lacking 'gold standard' reference SVs to help distinguish true

26 from false calls. Most studies rely on combining an ensemble of signals from different SV detection

27 methods, although this strategy does not reliably improve performance and can in some cases aggravate

28 false discovery[3]. Researchers therefore often apply independent experimental[8-9] or visualization methods[10]

29 to validate a subset of SV calls. Overall, there remains an unsatisfactory lack of consensus on how to

30 validate the quality of *de novo* SV datasets in most species[3].

31

32 Salmonids have the highest combined economic, ecological and scientific importance among all fish

33 lineages, and have consequently been subject to hundreds of genetics studies employing SNPs and other

34 molecular markers[11,12]. In common with most non-model fish species, the SV landscape remains extremely

35 poorly characterized in salmonids, apart from recent work informed by SNPs that revealed multi-megabase

36 inversions in rainbow trout (*Oncorhynchus mykiss* Walbaum) influencing migration[13,14], and a

37 chromosomal fusion under selection in Atlantic salmon[15], consistent with roles in adaptation. Salmonids

38 offer a unique system to characterize SVs due to an ancestral salmonid-specific autotetraploidization (i.e.

39 whole genome duplication, WGD) event (Ss4R), which occurred 80-100 Mya, following an earlier WGD

40 (300-350 Mya) in the teleost common ancestor[16,17,18]. WGD events may influence selection on SV retention

41 due to the functional redundancy linked to mass retention of duplicated genes, though this idea is yet to be

42 tested. In addition, salmonids have been farmed in aquaculture for a small number (<15) of generations[11],

43 and while the genetic architecture of such recent domestication has been investigated using SNPs[19], the role

44 played by SVs remains unexplored. Finally, the application of SVs in selective breeding of salmonids and

45 other commercial fishes remains untested. Clearly, the lack of SV data and analysis frameworks in

46 salmonids represents an important knowledge gap.

47

48 Here we provide an end-to-end workflow to detect, genotype, validate and annotate SVs using short-read

49 sequencing, removing false positives through efficient manual curation[10], allowing reliable SV discovery in

50 non-model species. Using this approach, we report a detailed investigation of the genomic landscape of

51 SVs in the iconic Atlantic salmon, inclusive of 492 genomes representing wild and farmed genetic

52 diversity, and populations of both European and North American descent.


53 **Results**

54

55 **Accurate SV discovery in Atlantic salmon**

56 We developed a workflow for SV discovery using paired-end short-read sequencing data aligned to the

57 unmasked ICSASG_V2 reference assembly[17], which can be run in Snakemake[20] (Supplementary Figure 1).

58 The probabilistic tool Lumpy[21] was used for SV detection, which simultaneously draws on multiple

59 evidence and SVtyper[22] was used for genotyping. As *de novo* SV detection using short-read data is prone to

60 false positives[3,21,23], we added steps to avoid SV calling in complex regions of the genome where false

61 positive rates were predicted to be particularly high (proven below). This included regions of ≥100x

62 coverage (>10 times higher than the global average of 8.1x coverage across 492 samples), shown elsewhere

63 to be overwhelmingly false calls[3], as well as gap regions in the ICSASG_V2 assembly. These complex

64 regions were most prevalent in chromosome arms where rediploidization was delayed after Ss4R,

65 characterized by high sequence similarity among duplicated regions[17] (Supplementary Figure 2).

66

67   Rather than using evidence from additional SV detection tools as a filter for true SV calls, a strategy shown

68   elsewhere to be potentially unreliable[3], we applied a curation approach to the entire filtered SV dataset

69   using SV-plaudit[10]. SV-plaudit is a scalable framework for the rapid production of thousands of SV images

70   via Amazon web services[10] (examples: Supplementary Figures 3-8). This approach allowed us to efficiently

71   retain high-confidence SV calls, while excluding low confidence or ambiguous calls, on the basis of

72   available visual evidence drawn from paired-end and split-read alignments, in addition to read depth[10,21].

73   The Atlantic salmon individuals (Supplementary Table l) produced on average 55,754 SV calls (median:

74   55,041, SD: 10,051) before filtering complex regions and SV-plaudit curation (Supplementary Table 2).

75   Across all individuals, 165,116 unique SVs were detected (size: 100bp to 2 million bp), which included an

76   outlier peak of deletion SVs in the 1,432-1,436 bp size range (Supplementary Data 1; Supplementary

77   Figure 9).

78

79   Using SV-plaudit on the full set of SV calls allowed us to retain only high-confidence calls, quantify the

80   impact of filtering complex regions, and estimate a false discovery rate (FDR). The overall estimated FDR

81   was 0.91 (149,491/165,116 of calls had low confidence), in line with the highest estimates in the

82   literature[2,3,7]. In complex regions, the FDR was 0.992 (47,268/47,636 calls had low confidence). In the

83   remaining chromosome-anchored assembly, the FDR was 0.85, validating the usefulness of removing

84   complex genomic regions. Sequencing depth was not a reliable indicator of FDR (Supplementary Figure

85   10).  A final high-quality set of 15,483 unique SV calls (14,017 deletions, 1,244 duplications, 242

86   inversions) and their genomic location is visualized in Fig. 1a and 1b. The average size for deletions was

87   1,532 bp (100 to 1,946,935 bp; SD: 23,070 bp) and for duplications 8,183 bp (102 to 80,1673 bp; SD:

88   25,589 bp) (Fig. 1c, d). For inversions, the average size was 121,935 bp (113 to 1,796,230 bp; SD: 278,698

89   bp) (Fig. 1e). The outlier peak at 1,432-1436 bp remained in the high-confidence deletions (Fig. 1c).

90

91   To validate our SV discovery workflow we estimated the true positive rate for SV presence/absence and

92   genotype calls using the high-confidence data retained after the SV-plaudit step. We sequenced PCR

93   amplicons for 876 independent SV calls representing 168 unique SVs (108 deletions, 46 duplications, 15

94   inversions) (Supplementary Figure 11) at ≥50x coverage on the MinION platform. Across all SV calls, the

95   true positive rate was 0.88 for SV presence/absence and 0.81 for SV plus genotype. For deletion calls, the

96   true positive rate was 0.93 for presence/absence (520/559 calls) and 0.85 (475/559 calls) for genotype. For

97   duplications, the true positive rate was 0.81 for presence/absence (186/230 calls) and 0.74 (170/230 calls)

98   for genotype. For inversion calls, the true positive rate was 0.78 for presence/absence (68/87 calls) and

99   0.75(65/87 calls) for genotype. Full results are shown in Supplementary Table 3 (with examples in

100  Supplementary Figures 12, 13 and 14). In summary, SV-plaudit curation vastly reduced the FDR to

101  maintain predominantly true SV calls (provided in Supplementary Data 2).

102

103  To further confirm data quality, we asked if the high confidence SVs genotypes capture expected

104  population genetic structure (Fig. 1f-j). SV genotypes were used in principal component analyses (PCA) for

105  the different SV types (Fig. 1f-i). For all SV types, PC1 separated European and Canadian salmon,

106  consistent with past work e.g.[24,25]. Deletions achieved a better resolution for the sampled European

107  populations, with PC2 separating populations from Europe into distinct groups explained by latitude with

108  evidence of intermixing at middle latitudes in Norway (Supplementary Figure 15), as reported elsewhere[24].

109  All farmed salmon clustered with the wild populations from which they are descended. Farmed salmon

110  from Europe, including 13 farmed fish from Chile, clustered with wild salmon from Southern Norway,

111  while 7 Chilean farmed salmon clustered with Canadian salmon (Fig. 1c). Using the high-confidence

112  deletion genotypes, an admixture analysis was performed, which was consistent with the PC analysis (Fig.

113  1j). For comparison, we also performed PCAs using the raw unfiltered SV calls, plus the reduced subset

114  filtered for complex regions, which failed to capture the same population structure (Supplementary Figure

115  16). In summary, our final set of deletion genotypes capture expected population genetic structure at the

116  highest resolution. It is unclear if the weaker signal for duplications and inversions is linked to specific

117  properties of these markers, their comparatively lower number, or slightly lower genotyping accuracy.

118

119  **Annotation of Atlantic salmon SVs**

120  We used SnpEff[26] to annotate all high confidence SV calls against features in the ICSASG_v2 annotation.

121  Many SVs were located in intergenic and intronic regions (Supplementary Figure 17), with 62%, 3% and

122  2.5% within 5 kb of a protein-coding gene, long non-coding RNA gene or pseudogene, respectively.

123 Around half (49%) of all SVs overlapped one or more RefSeq gene, the majority of which overlapped a

124 single gene (Supplementary Figure 18), with 8,439 genes overlapped in total. Approximately 4%, 21% and

125 25% of deletions, duplications and inversions were predicted by SnpEff to have a high impact, respectively,

126 including hundreds of putative exon losses, frameshift variants and potential gene fusion events

127 (Supplementary Figure 19). 101 duplications spanned entire genes (mean length: 51.7 kb, median length:

128 15.1 kb). The high impact annotations for different SV types were associated with an overrepresentation of

129 several biological processes in the gene ontology (GO) framework[27] (Supplementary Table 4, 5).

130
131 **Recently active DNA transposon in *Salmo* evolution**

132 The outlier peak observed in the deletion calls (Fig. 1a; Supplementary Figure 9) was investigated by

133 extracting all high confidence variants of 1,432-1,436 bp in size (104 sequences) from the ICSASG_v2

134 genome. 94 and 89 of these sequences shared ≥50% and ≥95% identity in all pairwise combinations,

135 respectively. The 94 sequences were used as queries in BLASTn searches revealing that 91% (86 out of 94)

136 shared ≥95% identity to a pTSsa2 piggyBac-like DNA transposon (NCBI accession: EF685967)[28]. The

137 breakpoints in the outlier deletions SV match to the complete pTSsa2 sequence (Supplementary Data 3),

138 missing no more than a few bp at the 5' or 3' end. Consequently, the outlier deletion peak (Fig. 1a) appears

139 to largely represent an intact pTSsa2 sequence.

140

141 Phylogenetic analysis was done incorporating the Atlantic salmon pTSsa2 sequences along with the top

142 100 BLASTn hits to the pTSsa2 sequence in the genome of brown trout *Salmo trutta* (repeat masking off;

143 all sequences e-value = 0.0, 70-100% and 84-95% query coverage and identity, respectively).  Repeating

144 the search against genomes for the next most closely-related salmonid genera, *Salvelinus* (Arctic charr *S.*

145 *alpinus*) and *Oncorhynchus* (rainbow trout *O. mykiss*, coho salmon *O. kitsuch,* and chinook salmon *O.*

146 *tshawytscha*) failed to identify sequences sharing >50% coverage or >81% identity. The tree indicates

147 independent expansions of pTSsa2 sequences in the Atlantic salmon and brown trout genome (Fig. 2;

148 Supplementary Figure 20). The pTSsa2 sequence appears in the Atlantic salmon genome with high copy

149 number across all chromosomes (Supplementary Figure 21).

150

151 We also determined the broader overlap of SVs and repeat sequences in the Atlantic salmon genome.

152 Among all SVs, 65% (10,184) contained no repeat sequences, 16% (2,423) a single repeat, and 7% (1,027)

153 two repeats. There was a significant correlation between SV size and the number of repeats per SV across

154 all SV types (Pearson's R $\geq$0.99, $P < 0.0001$ in each test), indicating that the number of repeats within each

155 SV was simply a direct product of SV size.

156

157 **Impact of genome duplication on the SV landscape**

158 Salmonid genomes retain a global signature of duplication from Ss4R, with at least half of the protein-

159 coding genes retained as expressed, functional duplicates (referred to as ohnologs)[17,18]. Ss4R ohnolog pairs

160 share amino acid sequence identity ranging from ~75 to 100%[12,17,18] with ~40% maintaining the ancestral

161 tissue expression pattern[17], suggesting pervasive functional redundancy. We hypothesised that the

162 redundancy provided ohnolog retention after WGD influenced the evolution of the SV landscape by

163 creating a mutational buffer[29] against deleterious SV mutations. A key prediction is that genes found in

164 Ss4R ohnolog pairs (with scope for functional redundancy) should be more overlapped by SVs compared

165 to singleton genes (lacking scope for functional redundancy).

166

167 We tested this prediction by generating a novel set of high-confidence Ss4R ohnolog pairs (10,023 pairs,

168 i.e. 20,046 genes) and singletons (8,282 genes) (Supplementary Data 4) and indeed, found a significant

169 enrichment of SVs overlapping retained Ss4R ohnologs (*Fisher's exact test*, $P = 0$, odds ratio = 1.47)

170 (Supplementary Table 6). This effect was specific to deletions ($P = 0$, odds ratio = 1.62), and hence not

171 observed in duplications ($P = 0.62$) nor inversions ($P = 0.52$). SVs with putative high impact did not

172 overlap ohnologs more than singletons (high impact snpEff annotation: $P = 0.93$, manually curated

173 deletions impacting exons: $P = 0.55$) (Supplementary Data 5).

174

175 Next we asked if gene expression characteristics influence the overlap between SVs and Ss4R ohnologs.

176 We initially used Spearman's rank correlation to establish co-expression of ohnologs across an RNA-Seq

177 atlas of 15 tissues[16]. We found that ohnolog pairs where one copy overlaps an SV showed slightly lower

178 expression correlation compared to randomly selected ohnolog pairs (resampling test, $P = 0$)

9

179 (Supplementary Figure 22). This pattern could be explained by SVs affecting ohnolog pairs with greater

180 levels of functional divergence, but may also be caused by relaxed purifying selection on duplicated copies,

181 allowing more SVs to accumulate. It has been shown elsewhere that the more highly expressed ohnolog in

182 a pair is typically under stronger purifying selection[30]. Therefore, we asked if ohnologs overlapped by a

183 deletion SV have reduced expression compared to their duplicate with no SV overlap. Indeed, this was the

184 case (*Wilcoxon rank-sum test*, $P$ = 2.9e-6) (Supplementary Figure 22). We also found that ohnolog pairs

185 showing overlap with deletion SVs showed reduced expression compared to ohnolog pairs showing no

186 overlap to SVs (*Wilcoxon rank-sum* test, $P$ = 7e-25) (Supplementary Figure 22).

187

188 Overall, these analyses reveal that the Ss4R WGD strongly influenced the retention of deletion SVs in the

189 Atlantic salmon genome, and this may be explained by functional redundancy.

190

191 **Selection on SVs during Atlantic salmon domestication**

192 Our study provides a unique opportunity to ask if SVs were selected during the domestication of Atlantic

193 salmon, which commenced when the Norwegian aquaculture industry was founded in the late 1960s[11,31].

194 Consequently, farmed Atlantic salmon are no more than 15 generations 'from the wild', in contrast to

195 livestock and poultry, which have been domesticated for thousands of years[11,12]. The early domestication

196 process involves strong selection on behavioural traits[32,33] targeting molecular pathways underpinning

197 cognition, learning and memory, for instance genes with functions in synaptic transmission and

198 plasticity[34,35]. Specifically, selection on farmed animals should remove individuals that invest in costly

199 behavioural and stress responses such as predator avoidance and fear processing, in favour of animals that

200 invest into performance traits [32,36]. We thus hypothesised that SVs linked to genes regulating pathways

201 controlling behaviour would be under distinct selective pressures in farmed and wild salmon.

202

203 To test our hypothesis, we established significantly genetically differentiated SVs by calculating the

204 fixation index ($F_{ST}$)[37] between 34 farmed Norwegian salmon and 257 wild salmon from Norway. The wild

205 individuals were selected based on a PCA including all European salmon, aiming to remove confounding

206 effects of genetic differentiation by latitude observed in wild Norwegian salmon (Fig. 3a), retaining the

10

207  closest possible background to the wild founders used in aquaculture. We used a permutation approach to

208  estimate the probability of observed $F_{ST}$ values in relation to random expectations, defining 584 SV outliers

209  at $P<0.01$ (all $F_{ST}$ >0.103, Median $F_{ST}$ = 0.149) (Fig. 3b; Supplementary Data 6), which were distributed

210  throughout the genome (Fig. 3c).

211

212  GO enrichment tests identified 132 overrepresented biological processes ($P<0.05$) among the genes linked

213  to these outlier SVs by SnpEff (Supplementary Table 7). This set comprises 326 unique genes contributing

214  to the enriched terms (Supplementary Table 8). 34 biological processes explained by 156 unique genes

215  (48% of the unique genes contributing to all enriched GO terms) were daughter terms related either to

216  learning and behaviour, including 'habituation' ($P<0.002$), 'vocal learning' ($P<0.001$), and 'adult behavior'

217  ($P<0.02$), or the nervous system, including 'positive regulation of nervous system process' ($P<0.02$),'

218  presynaptic membrane assembly'($P<0.01$), 'postsynapse assembly' ($P<0.02$) 'oligodendrocyte

219  development' ($P<0.001$) and 'regulation of neuronal synaptic plasticity' ($P<0.03$).

220

221  To test our hypothesis, we asked if genes linked to outlier SVs showed enrichment in brain expression (Fig.

222  3d). Indeed, this was strongly supported when judged against transcriptome-wide expectations (Fig. 3d);

223  with the signal being strongest for the 326 gene subset contributing to the overrepresented GO terms,

224  emphasising particular importance of brain functions among the enriched gene set (Fig. 3d, Supplementary

225  Table 9). A positive enrichment in the expression of outlier linked genes was only observed in brain, with

226  nine other tested tissues showing either no differences to transcriptomic expectations, or in the case of

227  muscle and foregut, reduced expression specificity (Supplementary Table 9; Supplementary Figures 23,

228  24). Finally, we asked if the outlier SVs overlapped putative cis-regulatory elements (CREs) detected in

229  brain using novel ATAC-Seq data (significant peaks overlapping a gene +/- 3,000bp up/downstream; n=4)

230  more than expected. For 9,920 SVs lacking evidence for differentiation between farmed and wild fish ($F_{ST}$

231  $P$ >0.05), 7.1% overlapped at least one brain ATAC-Seq peak, which was almost identical to SV outliers

232  (7.0%) (*Fisher's exact test*, $P = 0.86$). A similar result was observed by restricting the analysis to genes

233  with brain biased expression (*Fisher's exact test*, $P = 0.41$).

234

11

235 **SVs selected by domestication are linked to many synaptic genes**

236 The increased brain expression and overrepresentation of nervous system functions for SV outlier linked

237 genes motivated us to investigate the role of these loci in the genetic architecture of domestication. We

238 performed a detailed annotation of the 156 SV outlier linked genes contributing to the 34 aforementioned

239 enriched GO terms (Supplementary Table 10). To cement the relevance of this gene set to our hypothesis,

240 we cross-referenced all the encoded protein products with a high-resolution synaptic proteome from

241 zebrafish[38]. Our rationale was that the synaptic proteome is central to nervous system activity and defines

242 the repertoire of cognitive and behaviours an animal can perform during its life[38,39].

243

244 Among the 156 SV outlier linked genes, 65 (i.e. 42%, linked to 67 distinct SVs) encode a protein with an

245 ortholog in the zebrafish synaptic proteome (Supplementary Table 10) defined by stringent reciprocal

246 BLAST (mean respective pairwise % identity and coverage = 77% and 95%). As synaptic proteomes are

247 highly conserved between fish and mammals[38], it is reasonable to assume these proteins are *bone fide*

248 components of Atlantic salmon synaptic proteomes, and that a minimum of 11% of the outlier SVs were

249 linked to synaptic genes by SnpEff. These proteins are encoded by multiple members of ancient, conserved

250 gene families involved in synaptic formation, transmission and plasticity, including neurexins (*NRXN1* and

251 *NRXN2*), SH3 and multiple ankyrin repeat domains 3 proteins (*SHANK2* and *3*), cadherins (*CDH4*, *CDH8*,

252 *CDH11*, *PCDH1*), Down syndrome cell adhesion molecules (*DSCAM* and *DSCAML*), teneurins (*TENM1*

253 and *TENM2*), gamma-aminobutyric acid receptors (*GABRB2* and *GABRG2*), potassium voltage-gated

254 channel subfamily D members (*KCND1* and *KCND2*), receptor-type tyrosine-protein phosphatases

255 (*PTPRG* and *PTPRN2*) and ionotropic glutamate receptors (*GRIK3* and *GRIN2C*) (Fig. 4). Genetic

256 disruption to orthologs for most these proteins (59/65) cause behavioural and/or neurological disorders in

257 mammals (Supplementary Table 10).

258

259 To ask how selection acted on these variants during domestication, we compared allele frequencies

260 between wild and farmed fish (Fig. 4). By far the most common scenario was that the synapse gene - linked

261 SVs are rare alleles in wild fish that show increased frequency of heterozygotes (carrying one SV copy,

262 0/1) and homozygotes (carrying both SV copies, 1/1) in farmed fish (Fig. 4). We also found that farmed

12

263  individuals often carry multiple copies of SVs that are especially rare in wild fish (defined as 0/0

264  homozygous frequency ≥ 0.90, 45 SVs) - assumed to be deleterious in natural environments - including

265  homozygote 1/1 states for SVs located on different chromosomes (Supplementary Figure 25).

266

267  Many of the outlier SVs linked to the 65 synaptic genes are located in non-coding regions (introns and

268  untranslated regions, 45%), while a smaller fraction are located within 10kb up or downstream (15%) or

269  within ≥10kb to 260 kb (33%) of the same genes (Fig. 4). A smaller fraction affect coding regions via

270  whole gene duplications, either involving a small number of genes, e.g. a 55 kb duplication overlapping the

271  brain-specific *CDK5R1* gene, or through larger multigene duplications (Fig. 4; Supplementary Table 10). A

272  striking example of an SV with a putative major disruptive effect was a 696 kb inversion that flips multiple

273  exons and the upstream region of the brain-specific gene encoding neurexin 2, which should halt translation

274  of a functional protein (Supplementary Table 10). Finally, among this synaptic gene set, we identified two

275  ohnolog pairs retained from Ss4R encoding astrotactin-1 and seizure protein 6 (Fig. 4).

276

277  **Major effect SVs altered by domestication**

278  We identified 32 further SVs with major predicted effects on gene structure and function among the

279  significant $F_{ST}$ outliers, which typically show increased allele frequency in farmed compared to wild Atlantic

280  salmon (Table 1). These SVs disrupt or ablate coding genes with diverse functions, including male fertility

281  (e.g. *CATSPERB*[40]), immunity (e.g. B cell survival and signalling, *GIMAP8*[41] and two distinct *CD22*[42] genes),

282  circadian control of metabolism (*NR1D2*[43]), lipid metabolism and insulin sensitivity (*ELOVL6*[44]), and

283  melanin transport and deposition (*MYRAP*[45]) (Table 1). We observed four deletions that disrupt conserved

284  lncRNAs of unknown function, and several large SVs that cover multiple genes, for instance a 423 kb

285  inversion on Chromosome 7 containing 16 genes that was absent in 257 wild salmon (Table 1). In summary,

286  this data demonstrates that diverse gene functions beyond neurological and behavioural pathways were

287  altered by the domestication of Atlantic salmon due to altered selective pressure or drift.

288 **Discussion**

289

290  Despite an increasing shift towards the use of long-read sequencing for SV discovery[1,2], these technologies

291  remain prohibitively expensive for large-scale population genetics, making such datasets scarce in most

292  species. Consequently, it remains a timely challenge to extract reliable SV calls from the more extensive

293  repository of short-read genome sequencing datasets, which continue to emerge rapidly in many species,

294  largely for use in SNP analyses. The approach reported can be applied for reliable SV detection and

295  genotyping using such data in any species with a reference genome. A critical step - unique to this study -

296  was the curation of all SV calls using SV-plaudit[10]. This approach demands significant manual effort,

297  equivalent to approximately two weeks for a small team of trained curators, yet was efficient in retaining

298  predominantly true calls, and allowed us to demonstrate the value of filtering complex regions to drastically

299  reduce the FDR. The overall extreme FDR for SV discovery advocates for the routine application of such

300  curation in SV studies based on short-read sequencing, particularly if 'gold-standard' SVs defined by past

301  work are unavailable.

302

303  The SVs reported provide a novel resource for future studies on the genetic architecture of traits in Atlantic

304  salmon, which has excluded SVs until now. It will be useful to overlap our SVs with genomic regions of

305  interest such as QTLs defined by SNPs, to investigate SVs as putative causal variants. For example, we

306  discovered a duplication on chromosome 14 that likely destroys the function of the *MYRIP* gene, which is

307  involved in melanosome transport[45] – a past study discovered a single QTL on chromosome 14 that

308  explained differences in melanocyte pigmentation between wild and domesticated fish[46], which may be

309  linked to this newly discovered SV. It will also be useful in future studies to apply SV markers directly in

310  genome wide association analyses, and to test their value for genomic prediction in salmon breeding

311  programmes[11,12]. While our study captured hundreds of Atlantic salmon genomes representing several

312  major phylogeographic groups, it fails to capture broader genetic diversity within this species, and due to

313  the retention of only high confidence SV calls, our method may be prone to false negatives. Further,

314  inherent limitations of short read sequencing data for SV detection presumably obscures detection of many

14

315     SVs, suggesting future SV studies in Atlantic salmon must also focus on adapting long-read sequence data,

316     and integrating short and long-read data for optimal SV discovery[1].

317

318     We discovered intact pTSsa2 polymorphisms within our SV dataset, and provided evidence for transposon

319     expansion after the split of *S. salar* and *trutta* ~10 Mya[16] (Fig. 2). The pTSsa2 transposon appears with high

320     copy number in the Atlantic salmon genome, suggesting an important role in shaping very recent genome

321     architecture. Transposons have largely been excluded from studies of contemporary genetic variation in

322     salmonids, but were central to genome rediploidization after the Ss4R WGD[17], and likely contributed to the

323     evolution of the sex determining locus, e.g.[47]. As work in other taxa has revealed that transposon

324     polymorphisms contribute to adaptive evolution [48,49] and speciation[50], future studies on pTSsa2 should

325     investigate such possibilities in *Salmo*. We also showed that Atlantic salmon deletion SVs are more likely

326     to overlap genes retained as ohnolog pairs from the Ss4R WGD event compared to singleton genes. This

327     supports the hypothesis that WGD events buffer against potential deleterious impacts of SVs on gene

328     function and regulation, consistent with past work[29,51]. However, the link between SVs and the Ss4R WGD

329     requires further investigation to more fully dissect the role of selection and drift in driving SV retention.

330

331     We discovered many SVs showing genetic divergence between farmed and wild Atlantic salmon linked to

332     synaptic genes responsible for behavioural variation[38,39]. Most were rare alleles in wild fish and showed a

333     small to moderate increase in frequency in domesticated populations, consistent with a polygenic genetic

334     architecture for behavioural traits altered by domestication, including risk-taking behaviour, aggression,

335     and boldness[32,52,53,54,55,56], affecting many unique genes from the same functional networks, mirroring the

336     polygenic basis for many human neurological traits[57,58,59]. The disruption of mammalian orthologs for many

337     of the same synaptic genes cause disorders including schizophrenia, intellectual disability, autism, and

338     Alzheimer's (Supplementary Table 10). For Atlantic salmon, we did not establish if these SVs are

339     causative variants or in linkage disequilibrium with other variants under selection. In several cases, it is

340     likely that the SVs discovered are causative variants due to their disruptive nature on protein coding gene

341     sequence potential (e.g. Table 1), including the ablation of the key synaptic protein neurexin-2, which

342     caused autism-related behaviours when induced experimentally in mice[60]. However, as many of the outlier

15

343  SVs were located in non-coding regions, this points to regulatory effects on gene expression, which may

344  have minor or additive effects on behavioural traits. Future work should test whether the outlier SVs alter

345  the expression or function of synaptic genes and directly influence behavioural phenotypes. Beyond

346  neurological systems, domestication altered the frequencies of numerous major effect SVs disrupting genes

347  with diverse functional roles (Table 1), providing candidate causative variants for ongoing investigations

348  into diverse traits. For instance, an increased frequency of SVs ablating the *ELOVL6* and *NR1D2* genes in

349  domesticated fish, which play key roles in lipid metabolism, insulin resistance, and the coordination of

350  metabolic functions with the circadian clock[44,45], is highly consistent with a recent transcriptomic study

351  demonstrating altered metabolism linked to disrupted circadian regulation in domesticated compared to

352  wild Atlantic salmon[61].

353

354  To conclude, given the rapidly growing recognition of the importance of establishing the role of SVs in

355  adaptation and other evolutionary processes in natural populations[62,63], in addition to commercial variation

356  relevant to breeding of farmed animals[64,65], we anticipate that this reliable description of the SV landscape

357  in Atlantic salmon will encourage more studies exploiting SV markers to address both fundamental and

358  applied questions in the genetics of non-model species.

359

360  **Methods**

361

362  **Sequencing data**

363  Paired-end whole genome sequencing data (mean 8.1x coverage, 2 x 100-150 bp) was generated for 472

364  Atlantic salmon on several different platforms (Supplementary Table 1). DNA extraction, quality control

365  and sequencing library preparation followed standard methods. Wild Atlantic salmon were sampled either

366  during organized fishing expeditions or by anglers during the sport fishing season with DNA extracted

367  from scales. We sampled n=80 wild Canadian individuals from 8 sites, n=359 Norwegian individuals from

368  52 sites (including n=5 landlocked dwarf salmon), n=8 Baltic individuals from a single site and n=4 White

369  sea individuals from a single site. Whole genome sequencing data was generated for 21 farmed individuals

370  (n=12 individuals from Mowi ASA; n=9 samples from Xelect Ltd) and downloaded for a further 20

371  individuals (NCBI accession: PRJNA287458).

372

**SV detection and genotyping**

374  Sequence alignment to the unmasked ICSASG_V2 assembly (GCA_000233375.4)[17] was done using BWA

375  v0.7.13[66]. Reads were mapped to the complete reference, including unplaced scaffolds, with random

376  placement of multi-mapping reads[67]. Reads mapping to unplaced scaffolds were discarded. Alignments

377  were converted to BAM format in Samtools v0.1.19[68]. Alignment quality, batch effects and sample error

378  were further assessed using Indexcov goleft v0.2.1[69]. Gap regions were extracted and converted to BED

379  format using a Python script (Supplementary Note 1); SV calls overlapping these regions were identified

380  using Bedtools[70] and removed.  Sample coverage was estimated using mosdepth v0.2.3[71]. High depth

381  regions were defined as ≥100x coverage and removed; this cut-off was a compromise to avoid generating

382  too many false SV calls, balanced against the risk of losing real SVs.  High depth regions located within

383  100bp were merged. SV detection was done using the Lumpy-based tool Smoove V.2.3[21] with genotypes

384  called by SVtyper[22]. Gap and high-depth regions were combined into a single BED file, which can

385  optionally be used to exclude these locations from SV detection in Lumpy (-exclude option). All of the

386  above steps were combined in a Snakemake (v.3.11.0)[20] workflow, with the input being paired-end

387  sequencing data (FASTQ format), and the output a VCF file with SV locations and genotypes for all

388  individuals in a study (Supplementary Figure 1; Supplementary Note 2 provides Snakefile).

389

**SV-plaudit curation**

391  All 165,116 SV calls generated in the study were curated using SV-plaudit[10]. A plotCritic website was

392  setup on Amazon Web Services where variant images produced in samplot v1.01

393  (https://github.com/ryanlayer/samplot) were deployed. SV curation involved the random visualization of

394  one homozygous wild-type (0/0; lacking SV, identical to reference genome), two heterozygous (0/1, with

395  one SV copy) and two homozygous-alternate (1/1, with two SV copies) individuals per SV, done using

396  cyvcf2 v0.11.5[72]. With each image the question "is this variant real?" was answered (options: 'No', 'Yes',

397  or 'Maybe'). Only high confidence variants ('Yes') were kept for downstream analysis. Three different co-

17

398  authors (ACB, MKG, EP) team-curated the full SV set. 1,000 random plots were commonly curated by

399  each researcher to establish congruence in decision making, and there was 100% agreement concerning

400  high confidence ('Yes') variants. Subsequently the SV plots were divided randomly and each set validated

401  independently across the 3 researchers and then merged.

402

403  **SV annotation**

404  High confidence SVs retained following SV-plaudit curation were filtered to remove redundant SVs using

405  the Bedtools *intersect* function (90% reciprocal overlap), removing 133 SVs and leaving 15,483 SVs used

406  in further analysis (Supplementary Data 2). The association between SVs and RefSeq genes within the

407  ICSASG_v2 assembly was done using SnpEff[26] (default parameters). GO enrichment tests were done using

408  the 'weight01' algorithm and Fisher's test statistic in the TopGo package[73]. The background set was all

409  genes in the RefSeq annotation. The R package 'Ssa.RefSeq.db'

410  (https://gitlab.com/cigene/R/Ssa.RefSeq.db)[74] was used to retrieve GO annotations from the ICSASG_v2

411  genome. The overlap between SV locations and repeats in the ICSASG_v2 annotation was done using

412  Bedtools[61] against an existing database[17].

413

414  **Phylogenetic analyses**

415  pTSsa2 sequences including EF685967 were used in BLASTn[75] searches  against the NCBI nucleotide

416  database (restricted to Salmonidae) in addition to unmasked assemblies for Atlantic salmon (ISCASG_v2),

417  brown trout (GCA_901001165.1), Arctic charr (GCA_002910315.2), rainbow trout (GCA_002163495.1),

418  chinook salmon (GCA_002872995.1) and coho salmon (GCA_002021735.2). Sequence alignments were

419  performed using Mafft[76] with default settings. Phylogenetic analysis was done using the IQTREE server[77]

420  with estimation of the best-fitting nucleotide substitution model (Bayesian Information Criterion) and 1,000

421  ultrafast bootstraps[78].

422

423  **SV validation by MinION sequencing**

424  PCR primers are shown in Supplementary Table 3. PCRs were performed using LongAmp® Taq (New

425  England Biolabs) with 1 cycle of 94ºC for 30s, 30 cycles of 94ºC for 30s, 56°C for 60s and 65 ºC for

18

426    50s/kb, followed by a 10 min extension at 65ºC. Amplicons for different SVs in each fish individual were

427    pooled and cleaned using AMPure XP beads (Beckman Coulter). 250ng pooled DNA was used to create

428    sequencing libraries with a 1D SQK-LSK109 kit (Oxford Nanopore Technologies, ONT). DNA was end-

429    repaired using the NEBNext Ultra II End Repair/dA Tailing kit (New England Biolabs) and purified using

430    AMPure XP beads. Native barcodes were ligated to end-repaired DNA using Blunt/TA Ligation Master

431    Mix. Barcoded DNA was purified with AMPure XP beads and pooled in equimolar concentration to a total

432    of 200 ng per library (~0.2 pmol). AMII Adapter mix (ONT) was ligated to the DNA using Blunt/TA

433    Ligation Master Mix (New England Biolabs) before the adapter-ligated library was purified with AMPure

434    XP beads. DNA concentration was determined at each step using a Qubit fluorimeter (Thermo Fisher

435    Scientific) with a ds-DNA HS kit (Invitrogen).

436

437    Sequencing libraries were loaded onto MinION FLO-MIN106D R9.4.1 flow cells (ONT) and run via

438    MinKNOW for 36h without real-time basecalling. Basecalling and demultiplexing was performed with

439    Guppy v2.3.7. FASTQ files were uploaded into Geneious Prime 2019.1.1 and simultaneously mapped to a

440    reference of sequences spanning all candidate SV regions in the ISCSAG_v2 assembly. Mapping was done

441    with the following parameters: 'medium-fast sensitivity', 'finding structural variants', including 'short

442    insertions' and 'deletions' of any size, with the setting 'map multiple best matches' set to 'None', and the

443    minimum support for SV discovery set to 2 reads. Alignments were inspected for the presence and

444    genotype of the SV. Amplicons with <50x coverage to the target SV region were discarded as failed PCRs.

445    When alignments matched the predicted SV breakpoints and size, the SV call was considered correct.

446    When >90% of the aligned reads matched to the expected SV and breakpoints (i.e. a gap for deletions, an

447    insertion for duplications and flipped reads for inversions compared to the reference) it was classified 1/1

448    homozygous. When at least 10% of the aligned reads matched to both the reference genome state, in

449    addition to the 1/1 state, the locus was classified 0/1 heterozygous.

450

451    **Association between SVs and Ss4R ohnologs**

452    The code used to identify a genome-wide set of Ss4R ohnologs, along with a description of the genome

453    assembly annotations employed, is available at https://gitlab.com/sandve-lab/salmonid_synteny and

454    https://gitlab.com/sandve-lab/defining_duplicates. Orthogroups were constructed with Orthofinder[79] using

455    seven salmonid species (Atlantic salmon, rainbow trout, Arctic charr, coho salmon, huchen *Hucho hucho*,

456    and European grayling *Thymallus thymallus*), five additional actinopterygians (zebrafish, medaka *Oryzias*

457    *latipes*, northern pike *Esox lucius*, three-spined stickleback *Gasterosteus aculeatus* and spotted gar

458    *Lepisosteus oculatus*), and two mammals (human and mouse *Mus musculus*). For each orthogroup, we

459    extracted nucleotide protein coding sequences, aligned them with Macse[80] and built gene trees using

460    TreeBeST[81]. Trees were split into smaller subtrees at the node representing the divergence between pike

461    and salmonids. To derive a final set of Atlantic salmon Ss4R ohnologs, we used both synteny and gene tree

462    topology criteria. Firstly, we required that the subtrees branched with northern pike as the sister to

463    salmonids and outgroup to Ss4R[16,17] and contained either exactly two (ohnologs) or exactly one (singletons)

464    Atlantic salmon genes. Secondly, we removed any putative Ss4R ohnologs falling outside conserved

465    synteny blocks predicted using iadhore[82]. A final set of ohnolog pairs is provided in Supplementary Data 4,

466    which contains all gene trees in NWK format.

467

468    We used the *fisher.exact()* function in R to compare the observed counts of SVs overlapping singleton and

469    ohnologs with the total counts of singletons and ohnologs. To test for association between ohnolog

470    expression divergence and SV overlap, we used a 15 tissue RNA-Seq dataset[17] available as a TPM

471    (transcripts per million reads) table in the salmofisher R-package https://gitlab.com/sandve-

472    lab/salmonfisher. We used the *cor()* function in R to compute median Spearman's tissue expression

473    correlation for all ohnolog pairs where one copy was overlapped by an SV. We then computed median

474    correlations for 1,000 randomly sampled ohnolog sets of the same size. The *P*-value was estimated as the

475    proportion of resampled medians lower than the observed median for ohnologs overlapped by SVs. Tests

476    comparing expression level between genes that were either overlapped or not overlapped by SVs were

477    conducted using the sum log10 transformed TPM for each gene across all 15 tissues. The function

478    *wilcox_test* within the R-package *rstatix* was used to calculate *P*-values for differences in expression levels.

479    The code used is available at https://gitlab.com/ssandve/atlantic_salmon_sv_ohnolog_analyses/.

480

481    **Association of SVs with brain ATAC peaks**

20

482   Four Atlantic salmon (freshwater stage, 26-28g) were killed using a Schedule 1 method following the

483   Animals (Scientific Procedures) Act 1986. Around 50mg homogenized brain tissue was processed to

484   extract nuclei using the Omni-ATAC protocol for frozen tissues[83]. Nuclei were counted on an automated

485   cell counter (TC20 BioRad, range 4-6 um) and further confirmed intact under microscope. 50,000 nuclei

486   were used in the transposition reaction including 2.5 µL Tn5 enzyme (Illumina Nextera DNA Flex Library

487   Prep Kit), incubated for 30 minutes at 37 °C in a shaker at 200 rpm. The samples were purified with the

488   MinElute PCR purification kit (Qiagen) and eluted in 12µL elution buffer. qPCR was used to determine the

489   optimal number of PCR cycles for library preparation[84] (8-10 cycles used). Sequencing libraries were

490   prepared with short fragments and fragments >1,000 bp removed using AMPure XP beads (Beckman

491   Coulter, Inc.). Fragment length distributions and confirmation of nucleosome banding patterns were

492   determined on a 2100 Bioanalyzer (Agilent) and the library concentration estimated using a Qubit system

493   (Thermo Scientific). Libraries were sent to the Norwegian Sequencing Centre, where paired-end 2 x 75 bp

494   sequencing was done on an Illumina HiSeq 4000. The raw sequencing data is available through

495   ArrayExpress (Accession: E-MTAB-9001).

496

497   ATAC-Seq reads were aligned to the Atlantic salmon genome (ICSASG_v2) using BWA (v0.7.17)[66] and a

498   merged peak set called combining the four replicates using Genrich (https://github.com/jsh58/Genrich)

499   with default parameters, apart from "-m 20 -j" (minimum mapping quality 20; ATAC-Seq mode). Bedtools

500   was used to identify SVs overlapping ATAC-Seq peaks (filtered at corrected $P \leq 0.01$) associated to genes,

501   defined as being located within 3,000 bp up/downstream of the start and end coordinates of the longest

502   transcript per gene.

503

504   **Population structure analyses and $F_{ST}$ analyses**

505   PCAs were performed separately on the complete set of high confidence deletions (14,017), duplications

506   (1,244) and inversions (242) using the *prcomp* and *autoplot* functions within GGplot2[85] in R. Genotypes

507   were coded into bi-allelic marker format to be compatible with standard population genetics methods.

508   Population structure was examined using NGSadmix[86] tested for group sizes of K=2-4.

509

510  $F_{ST}$ values were calculated for all high confidence SVs using VCFtools v0.1.16[87] with the Weir and

511  Cockerham method[37] comparing 34 Norwegian farmed vs. 257 Norwegian wild Atlantic salmon (Fig. 4a

512  provides rationale for sample selection). To establish the significance of each $F_{ST}$ value, individuals from

513  the two groups were randomly split into two sets of the original size (i.e. 34 vs. 257 individuals) 200 times,

514  before the distribution of resultant $F_{ST}$ values was plotted using the ggplot2 function *geom_freqpoly*

515  (binwidth = 0.01). Per SV *P*-values were considered as the proportion of $F_{ST}$ values obtained in the 200

516  random distributions higher than the $F_{ST}$ in the observed distribution. Thus, if 10/200 randomly sampled

517  $F_{ST}$ values above the observed $F_{ST}$ value were recorded, *P*=0.05 was assigned. We further applied an $F_{ST}$

518  cutoff to include SVs where 99.7% of all $F_{ST}$ values fell above the randomly sampled values ($F_{ST}$ >0.103).

519  Any SVs lacking alternative alleles in the compared groups were excluded. Code to perform these analyses

520  is provided in Supplementary Note 3.

521

522  **Annotation of SV outliers**

523  GO enrichment tests for genes linked to the SV outliers (*P* <0.05) were done as described in the section

524  'SV annotation', with the background gene set restricted to all RefSeq genes linked to SVs by SnpEff. To

525  investigate the expression of genes linked to SV outliers, we used existing RNA-seq data[17], representing

526  normalized counts per million (CPM) for 10 tissues (brain, liver, muscle, spleen, pancreas, heart, pyloric,

527  gill, skin and foregut). We filtered any genes where the across-tissue sum of CPM was <1.0. A 'tissue

528  specificity' index was calculated, representing the sum across-tissue CPM divided by the CPM per tissue.

529  We tested whether genes linked to SV outliers by SnpEff, in addition to a subset contributing to significant

530  GO terms (*P*<0.01), differed from the transcriptome-wide expectations. Hypergeometric tests were used

531  (*dhyper* function in R) to compare the number of genes in the two gene sets with a tissue specificity index

532  ≥0.5 compared to all genes in the transcriptome. Two-sample t-tests (*t.test* function in R) were used to

533  compare differences in mean CPM between the two gene sets compared to all genes in the transcriptome.

534  BLAST was used to cross-reference protein products of genes linked to SV outliers against 3,840 unique

535  proteins detected in the zebrafish synaptic proteome[38] (downloaded from the GRCz11 assembly version

536  using BioMart at Ensembl.org), taking forward the top zebrafish BLAST hit (cut-off: 40% identity, 40%

537  query coverage) as a query in a reciprocal BLAST against all *S. salar* RefSeq proteins (no cut-off);

22

538 evidence for orthology was accepted when the candidate zebrafish protein showed a best hit to the original

539 query in the complete salmon proteome. We used the *fisher.exact()* function in R to test if the 584

540 significant $F_{ST}$ outlier SVs were more likely to overlap brain ATAC-Seq peaks than non-significant SVs (*P*

541 > 0.05), which was done considering all expressed genes (TPM ≥1) in the RNA-Seq tissue atlas described

542 above[17] and a subset of the same genes most highly expressed in brain (filtered for genes where brain was

543 among the top 3 tissues for TPM). The bedtools[61] intersect function was used to associate ATAC-Seq peaks

544 with SVs. The code used is available at https://gitlab.com/ssandve/atlantic_salmon_sv_ohnolog_analyses/.

545

546 **Data availability**

547 New genome sequences generated are available through the European Nucleotide Archive (project

548 accession: PRJEB38061, released upon publication). Sample accession numbers for all 492 Atlantic salmon

549 genomes are provided in Supplementary Table 1 (available upon publication). ATAC-Seq reads were

550 deposited in ArrayExpress (accession: E-MTAB-9001).

551

552 **Code availability**

553 Python script used to identify regions in ICSASG_v2 genome and convert output to BED file:

554 Supplementary Note 1.

555 Snakefile and associated code for SV detection pipeline: Supplementary Note 2.

556 R script used to obtain $F_{ST}$ values from random comparisons and establish probability value for outlier SVs:

557 Supplementary Note 3.

558 Code to define orthogroups and build gene trees: https://gitlab.com/sandve-lab/salmonid_synteny.

559 Code to identify Atlantic salmon ohnolog pairs from ortholog groups and gene trees:

560 https://gitlab.com/sandve-lab/defining_duplicates.

561 Code to analyse overlaps between SVs, ohnologs and ATAC-Seq data:

562 https://gitlab.com/ssandve/atlantic_salmon_sv_ohnolog_analyses/.

563

564 **Acknowledgements**

23

576

577  **Contributions**

578  DJM, SL, IAJ and SRS conceived the study. ACB, RL and TN developed the SV detection workflow. ACB

579  performed downstream analyses with contributions from MKG, DR, DJM, TDM, SRS, and EP. DJM (lead

580  supervisor), TJA, IAJ and SAM supervised ACB. ACB and MDG performed MinION sequencing. SLL

581  and MHH performed ATAC-Seq. KH, HS., BF-L, JE, CRP and LB provided wild Atlantic salmon samples.

582  ACB, DJM and SRS drafted the text and figures. All authors commented on and approved the final

583  manuscript.

584

585  **Ethics declarations**

586  Competing interests: The authors declare no competing interests

587

588  **Supplementary information**

589  Supplementary Figure 1. Snakemake pipeline for end-to-end SV detection.

590  Supplementary Figure 2. Locations of complex regions in Atlantic salmon genome.

591  Supplementary Figure 3. Example of SV-plaudit image for a high confidence deletion SV call

24

592    Supplementary Figure 4. Example of SV-plaudit image for a false positive deletion SV call excluded

593    from further analyses.

594    Supplementary Figure 5. Example of SV-plaudit image for a high confidence duplication SV call

595    retained in further analyses

596    Supplementary Figure 6. Example of SV-plaudit image for a false positive duplication SV call excluded

597    from further analyses

598    Supplementary Figure 7. Example of SV-plaudit image for a high confidence inversion SV call retained

599    in further analyses

600    Supplementary Figure 8. Example of SV-plaudit image for a false positive inversion SV call excluded

601    from further analyses

602    Supplementary Figure 9. SV sizes before SV-plaudit curation

603    Supplementary Figure 10. Sequencing depth was not a strong predictor of the final number of high-

604    confidence SVs retained after SV-plaudit curation.

605    Supplementary Figure 11. Summary of 168 SV regions used for MinION amplicon sequencing to

606    validate Lumpy/SVtyper SV and genotype calls.

607    Supplementary Figure 12. Example of congruence between SV/genotype calls and data generated by

608    MinION amplicon sequencing.

609    Supplementary Figure 13. Example of congruence between SV/genotype calls and data generated by

610    MinION amplicon sequencing.

611    Supplementary Figure 14. Example of congruence between SV/genotype calls and data generated by

612    MinION amplicon sequencing.

613    Supplementary Figure 15.  PCAs showing the same data presented in Fig. 1g-i (main text), except

614    visualized according to latitude

615    Supplementary Figure 16.  PCA analyses done on SV genotype calls prior to SV-plaudit curation

616    Supplementary Figure 17. SV annotation by SnpEff

617    Supplementary Figure 18. Overlap between high confidence SVs and protein coding genes in the

618    ICSASG_v2 annotation.

619 Supplementary Figure 19. Number of high impact annotations per snpEff effect for high-confidence

620 Atlantic salmon SVs.

621 Supplementary Figure 20. Maximum likelihood tree presented in Fig. 2 including sample identifiers,

622 genomic locations of pTSsa2 sequences and bootstrap values.

623 Supplementary Figure 21. Circos plot showing the genomic locations of pTSsa2 sequences in the

624 Atlantic salmon genome

625 Supplementary Figure 22. Expression characteristics of ohnologs depending on SV overlap.

626 Supplementary Figure 23. Tissue expression levels comparing SV outliers with transcriptome wide

627 expectations for nine tissues

628 Supplementary Figure 24. Tissue specificity comparing SV outliers with transcriptome wide

629 expectations for nine tissues

630 Supplementary Figure 25. Heatmap showing individual SV genotypes for 45 SV outliers linked to

631 synapse genes.

632

633 Supplementary Note 1. Python script used to extract gap regions in the ICSASG_v2 genome and and

634 convert the outputs to a BED file

635 Supplementary Note 2. Snakefiles and associated code for SV calling pipeline

636 Supplementary Note 3: Custom R script used to obtain $F_{ST}$ values from random comparisons and establish

637 probability value for outlier SVs

638

639 Supplementary Table 1. Details of samples used in study

640 Supplementary Table 2. SV call statistics per individual across 492 Atlantic salmon samples following

641 different filtering steps

642 Supplementary Table 3. Validation of SV calls and genotypes using MinION sequencing

643 Supplementary Table 4. GO Biological Process enrichment analysis for genes affected by high impact

644 deletions, duplications and inversions

645 Supplementary Table 5. Genes contributing to significant GO terms for high impact SVs

26

646    Supplementary Table 6. Fishers Exact test results contrasting the overlap between SVs with singleton genes

647    vs. Ss4R ohnolog genes.

648    Supplementary Table 7. GO enrichment analysis for genes linked to SV outliers between wild and farmed

649    Atlantic salmon

650    Supplementary Table 8. Genes contributing to significant GO terms for genes linked to SV outliers.

651    Supplementary Table 9. Statistical tests of two expression characteristics (specificity and level) across a

652    panel of tissues for 327 SV outlier linked genes contributing to significantly enriched GO biological

653    processes in comparison to a transcriptome-wide set gene set.

654    Supplementary Table 10. Detailed annotation of prioritized SV outliers between farmed and wild Atlantic

655    salmon linked to genes with synaptic functions.

656

657    Supplementary Data 1: Full SV dataset and genotypes prior to SV-plaudit curation

658    Supplementary Data 2: High-confidence SVs retained after SV-plaudit curation, including individual

659    genotypes and SnpEff annotation

660    Supplementary Data 3: Alignment of SV deletions representing pTSsa2 piggyBac-like DNA transposons

661    (used to Generate Fig. 2)

662    Supplementary Data 4: High confidence annotation of Ss4R ohnolog and singletons in the Atlantic salmon

663    genome

664    Supplementary Data 5: Manually filtered SV deletions that alter protein-coding exons

665    Supplementary Data 6: Significant SV outliers between wild and farmed salmon from Norway

666
667    **References**
668

669    1. Ho, S.S., Urban, A.E., Mills, R.E. Structural variation in the sequencing era. *Nat. Rev. Genet.* doi:
670    10.1038/s41576-019-0180-9 (2019).
671    2. Mahmoud, M., et al. Structural variant calling: the long and the short of it. *Genome Biol.* **20**, 246 (2019).
672    3. Cameron, D.L. Di Stefano, L., Papenfuss A.T. Comprehensive evaluation and characterisation of short
673    read general-purpose structural variant calling software. *Nat. Commun.* **10**, 3240 (2019).
674    4. Frazer, K.A., et al. Human genetic variation and its contribution to complex traits. *Nat. Rev. Genet.* **10**,
675    241-51 (2009).
676    5. Conrad, D.F., Hurles, M.E. The population genetics of structural variation. *Nat. Genet.* 39, S30–S36
677    (2007).
678    6. Chiang, C. et al. The impact of structural variation on human gene expression. *Nat. Genet.* 49, 692–699
679    (2017).

680   7. Kosugi, S., et al. Comprehensive evaluation of structural variation detection algorithms for whole
681   genome sequencing. *Genome Biol*. **20**, 117 (2019).
682   8. Becker, T., et al. FusorSV: an algorithm for optimally combining data from multiple structural variation
683   detection methods. *Genome Biol*. 19, 38 (2018).
684   9. Alkan, C., Coe, B.P, Eichler E.E. Genome structural variation discovery and genotyping. *Nat*. *Rev*.
685   *Genet*. **12**, 363-376 (2011).
686   10. Belyeu, J.R., et al. SV-plaudit: A cloud-based framework for manually curating thousands of structural
687   variants. *GigaScience* **7**, giy064 (2018).
688   11. Houston, R.D., et al. Harnessing genomics to fast-track genetic improvement in aquaculture. *Nat. Rev.*
689   *Genet.* Doi: s41576-020-0227-y (2020)
690   12. Houston, R.D., Macqueen, D.J. Atlantic salmon (*Salmo salar* L.) genetics in the 21st century: taking
691   leaps forward in aquaculture and biological understanding. *Anim*. *Genet*. **50**, 3-14 (2019)
692   13. Pearse, D.E., et al. Sex-dependent dominance maintains migration supergene in rainbow trout. *Nat.*
693   *Ecol*. *Evol*. **3**, 1731-42 (2019).
694   14. Pearse, D.E., et al. Rapid parallel evolution of standing variation in a single, complex, genomic region
695   is associated with life history in steelhead/rainbow trout. *Proc*. *Biol*. *Sci*. **281**, 20140012.
696   15. Wellband, C. et al. Chromosomal fusion and life history-associated genomic variation contribute to
697   within-river local adaptation of Atlantic salmon. *Mol. Ecol.* **28**, 1439-1459.
698   16. Macqueen, DJ., Johnston, I.A. A well-constrained estimate for the timing of the salmonid whole
699   genome duplication reveals major decoupling from species diversification. *Proc*. *Biol*. *Sci*. 2014 **281**,
700   20132881 (2014)
701   17. Lien, S., et al. The Atlantic salmon genome provides insights into rediploidization. *Nature* **533**,200-5
702   (2016).
703   18. Berthelot, C. et al. The rainbow trout genome provides novel insights into evolution after whole-
704   genome duplication in vertebrates. *Nat*. *Commun*. 2014 **5**, 3657 (2014)
705   19. López, M.E., et al. Comparing genomic signatures of domestication in two Atlantic salmon (*Salmo*
706   *salar* L.) populations with different geographical origins. *Evol. Appl.* **12**, 137-156 (2019)
707   20. Köster, J., Rahmann, S. Snakemake--a scalable bioinformatics workflow engine. *Bioinformatics* **28**,
708   2520-2 (2012).
709   21. Layer, R.M. et al. 2014. LUMPY: a probabilistic framework for structural variant discovery. *Genome*
710   *Biol*. **15**, R84 (2014).
711   22. Chiang et al. SpeedSeq: Ultra-fast personal genome analysis and interpretation. *Nat. Methods* **12**, 966–
712   968 (2015).
713   23. Kronenberg, Z.N. et al. Wham: Identifying structural variants of biological consequence. *PLoS.*
714   *Comput. Biol.* **11**, e1004572 (2015).
715   24. Wennevik, V., et al. Population genetic analysis reveals a geographically limited transition zone
716   between two genetically distinct Atlantic salmon lineages in Norway. *Ecol. Evol.* **9**, 6901–21. (2019)
717   25. Rougemont, Q., Bernatchez, L. The demographic history of Atlantic salmon (*Salmo salar*) across its
718   distribution range reconstructed from approximate Bayesian computations. *Evolution* **72**, 1261-77 (2018)
719   26. Cingolani, P. et al. A program for annotating and predicting the effects of single nucleotide
720   polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly*
721   (*Austin*) **6**, 80-92 (2012).
722   27. Ashburner, M. et al. Gene ontology: tool for the unification of biology. *Nat*. *Genet*. **25**, 25-9 (2000).
723   28. de Boer, J.G., et al. Bursts and horizontal evolution of DNA transposons in the speciation of
724   pseudotetraploid salmonids. *BMC Genomics* **8**, 422 (2007)
725   29. Fares, M. The origins of mutational robustness. *Trends Genet*. **31**, 373-381 (2015)
726   30. Pophaly, S.D., Tellier, A. Population Level Purifying Selection and Gene Expression Shape
727   Subgenome Evolution in Maize. *Mol*. *Biol*. *Evol*. **32**, 3226-35 (2015)
728   31. Gjedrem, T., Gjøen, H.M., Gjerde B. Genetic origin of Norwegian farmed Atlantic salmon.
729   *Aquaculture* **98**, 41-50 (1991)
730   32. Pasquet, A. Effects of Domestication on Fish Behaviour. In Book 'Animal Domestication'
731   10.5772/intechopen.78752 (InTechOpen, 2018).
732   33. Jensen, P.  Behavior genetics and the domestication of animals. *Annu*. *Rev*. *Anim*. *Biosci*. **2**, 85-104
733   (2014)

734  34. O'Rourke, T., Boeckx, C. Glutamate receptors in domestication and modern human evolution.
735  *Neurosci*. *Biobehav*. *Rev*. **108**, 341-357 (2020)
736  35. Theofanopoulou, C. et al. Self-domestication in *Homo sapiens*: Insights from comparative genomics.
737  *PLoS One* **12**, e0185306 (2017)
738  36. Price, E. O. Behavioral development in animals undergoing domestication. *Appl*. *Anim*. *Behav*. *Sci*.
739  245-271 (1999)
740  37. Weir, B.S., Cockerham, C.C. Estimating F-statistics for the analysis of population structure. *Evolution*
741  **38**, 1358-1370 (1984)
742  38. Bayés, À., et al. Evolution of complexity in the zebrafish synapse proteome. *Nat*. *Commun*. 2017 **8**,
743  14613 (2017)
744  39. Emes, R.D., Grant, S.G. Evolution of synapse complexity and diversity. *Annu*. *Rev*. *Neurosci*. **35**, 111-
745  31 (2012)
746  40. Liu, J. et al. CatSperbeta, a novel transmembrane protein in the CatSper channel complex. *J*. *Biol*.
747  *Chem*. **282**, 18945-52 (2007)
748  41. Webb, L.M., et al. Generation and characterisation of mice deficient in the multi-GTPase domain
749  containing protein, GIMAP8. *PLoS One* **9**, e110294 (2014)
750  42. Clark, E.A. and Giltiay, N.V. CD22: A Regulator of Innate and Adaptive B Cell Responses and
751  Autoimmunity. *Front*. *Immunol*. **9**, 2235 (2018)
752  43. Bugge, A. et al. Rev-erbα and Rev-erbβ coordinately protect the circadian clock and normal metabolic
753  function. *Genes Dev*. **26**, 657-67 (2012)
754  44. Matsuzaka, T. et al. Crucial role of a long-chain fatty acid elongase, Elovl6, in obesity-induced insulin
755  resistance. *Nat*. *Med*. **13**,1193-202 (2007)
756  45. Wasmeier, C. et al. Melanosomes at a glance. *J*. *Cell Sci*. **121**, 3995-9 (2008)
757  46. Jørgensen, K.M. et al. Judging a salmon by its spots: environmental variation is the primary
758  determinant of spot patterns in Salmo salar. *BMC Ecol*. **18**, 14 (2018)
759  47. Faber-Hammond, J.J., Phillips, R.B., Brown, K.H. Comparative analysis of the shared sex-
760  determination region (SDR) among salmonid fishes. *Genome Biol*. *Evol*. 7, 1972–1987 (2015)
761  48. Schrader, L., et al. Transposable element islands facilitate adaptation to novel environments in an
762  invasive species. *Nat*. *Commun*. **5**, 5495 (2014)
763  49. Bourgeois, Y., Boissinot, S. On the population dynamics of junk: a review on the population genomics
764  of transposable elements. *Genes (Basel)* **10**, pii:E419 (2019)
765  50. Laporte M, et al. DNA methylation reprogramming, TEs derepression and postzygotic isolation of
766  nascent species. *Sci*. *Adv*. 5, eaaw1644 (2019)
767  51. Gu, Z., et al. Role of duplicate genes in genetic robustness against null mutations. *Nature* **421**, 63-6
768  (2003)
769  52. Fleming I.A., Einum, S. Experimental tests of genetic divergence of farmed from wild Atlantic salmon
770  due to domestication. *ICES J*. *Mar*. *Sci*. **54**, 1051-1063 (1997)
771  53. Biro, P.A, et al. Predators select against high growth rates and risk-taking behaviour in domestic trout
772  populations. *Proc*. *R*. *Soc*. *B*. **271**, 2233–2237 (2004)
773  54. Lucas, M.D. et al. Behavioral differences among rainbow trout clonal lines. *Behav*. *Genet*. **34**, 355-65.
774  (2004)
775  55. Berejikian, B.A., et al. Competitive differences between newly emerged offspring of captive-reared and
776  wild coho salmon. *Trans*. *Am*. *Fish*. *Soc*. **128**, 832-839 (1999)
777  56. Solberg, M.F., et al. Domestication leads to increased predation susceptibility. *Sci*. *Rep*. **10**, 1929
778  (2020)
779  57. McCarroll, S.A., Hyman, S.E. Progress in the genetics of polygenic brain disorders: significant new
780  challenges for neurobiology. *Neuron* **80**, 578-87 (2013)
781  58. Lee, J.L. et al. Gene discovery and polygenic prediction from a genome-wide association study of
782  educational attainment in 1.1 million individuals. *Nat*. *Genet*. **50**, 1112-1121 (2018)
783  59. Purcell, S.M. et al. Common polygenic variation contributes to risk of schizophrenia and bipolar
784  disorder. *Nature* **460**, 748-52 (2009)
785  60. Dachtler, J. et al. Deletion of α-neurexin II results in autism-related behaviors in mice. *Transl*.
786  *Psychiatry* **4**, e484 (2014)
787  61. Jin, Y. et al. Comparative transcriptomics reveals domestication-associated features of Atlantic salmon
788  lipid metabolism. *Mol Ecol*. In Press doi: 10.1111/mec.15446 (2020)

29

789  62. Mérot, C., et al. A roadmap for understanding the evolutionary significance of structural genomic
790  variation. *Trends Ecol. Evol*. Doi: https://doi.org/10.1016/j.tree.2020.03.002 (2020)
791  63. Wellenreuther, M. et al. Going beyond SNPs: The role of structural genomic variants in adaptive
792  evolution and species diversification. *Mol. Ecol*. **28**, 1203-1209.
793  64. Bickhart, D.M., Liu, G.E. The challenges and importance of structural variation detection in livestock.
794  *Front. Genet*. **5**, 37 (2014)
795  65. Low, Y.W., et al. Haplotype-resolved cattle genomes provide insights into structural variation and
796  adaptation. *Nature Commun*. **11**, 2071 (2020)
797  66. Li, H., Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
798  *Bioinformatics* **25**, 1754-60.
799  67. Treangen, T.J., Salzberg S.L. Repetitive DNA and next-generation sequencing: computational
800  challenges and solutions. *Nat. Rev. Genet*. **13**, 36-46 (2011)
801  68. Li, H., et al. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-9
802  (2009)
803  69. Pedersen, B.S., et al. Indexcov: fast coverage quality control for whole-genome sequencing.
804  *Gigascience* **6**, 1-6. (2017)
805  70. Quinlan, A.R., Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic features.
806  *Bioinformatics* **26**, 841-2 (2010)
807  71. Pedersen, B.S., Quinlan, A.R. Mosdepth: quick coverage calculation for genomes and exomes.
808  *Bioinformatics* **34**, 867-868 (2018)
809  72. Pedersen, B.S. and Quinlan, A.R. cyvcf2: fast, flexible variant analysis with Python. *Bioinformatics* **33**,
810  1867-1869.
811  73. Alexa, A., Rahnenfuhrer, J. topGO: Enrichment Analysis for Gene Ontology. R package version 2.38.1
812  (2019)
813  74. Robertson, F.M. Lineage-specific rediploidization is a mechanism to explain time-lags between genome
814  duplication and evolutionary diversification. *Genome Biol*. **18**, 111 (2011)
815  75. Altschul, S.F., et al. Basic local alignment search tool. *J. Mol. Biol*. **215**, 403-10 (1990)
816  76. Katoh, K., Rozewicki, J., Yamada, K.D. MAFFT online service: multiple sequence alignment,
817  interactive sequence choice and visualization. *Brief. Bioinform*. **20**, 1160-1166 (2019)
818  77. Trifinopoulos, J., et al. W-IQ-TREE: a fast online phylogenetic tool for maximum likelihood analysis.
819  *Nucleic Acids Res*. **44**, W232-5 (2016)
820  78. Minh, B.Q., Nguyen, M.A., von Haeseler, A. Ultrafast approximation for phylogenetic bootstrap. *Mol.
821  Biol. Evol*. **30**, 1188-95.
822  79. Emms, D.M., Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative genomics.
823  *Genome Biol*. **20**, 238 (2019)
824  80. Ranwez, V. et al. MACSE: Multiple Alignment of Coding SEquences accounting for frameshifts and
825  stop codons. *PLoS One*. **6**, e22594 (2011)
826  81. Vilella, A.J. et al. EnsemblCompara GeneTrees: Complete, duplication-aware phylogenetic trees in
827  vertebrates. *Genome Res*. **19**, 327-35.
828  82. Proost, S. et al. i-ADHoRe 3.0--fast and sensitive detection of genomic homology in extremely large
829  data sets. *Nucleic Acids Res*. 40, e11 (2012)
830  83. Corces, M.R., et al. An improved ATAC-seq protocol reduces background and enables interrogation of
831  frozen tissues. *Nat. Methods* **14**, 959-62 (2017)
832  84. Buenrostro, J.D., et al. ATAC-seq: a method for assaying chromatin accessibility genome-wide. *Curr.
833  Protoc. Mol. Biol*. **109**, 21.29.1–21.29.9 (2015)
834  85. Wickham, H. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York. ISBN 978-3-
835  319-24277-4, https://ggplot2.tidyverse.org (2016)
836  86. Skotte, L., Korneliussen, T.S., Albrechtsen, A. Estimating individual admixture proportions from next
837  generation sequencing data. *Genetics* **195**, 693-702 (2013)
838  87. Danecek, P., et al. The variant call format and VCFtools. *Bioinformatics* **27**, 2156-8 (2011)
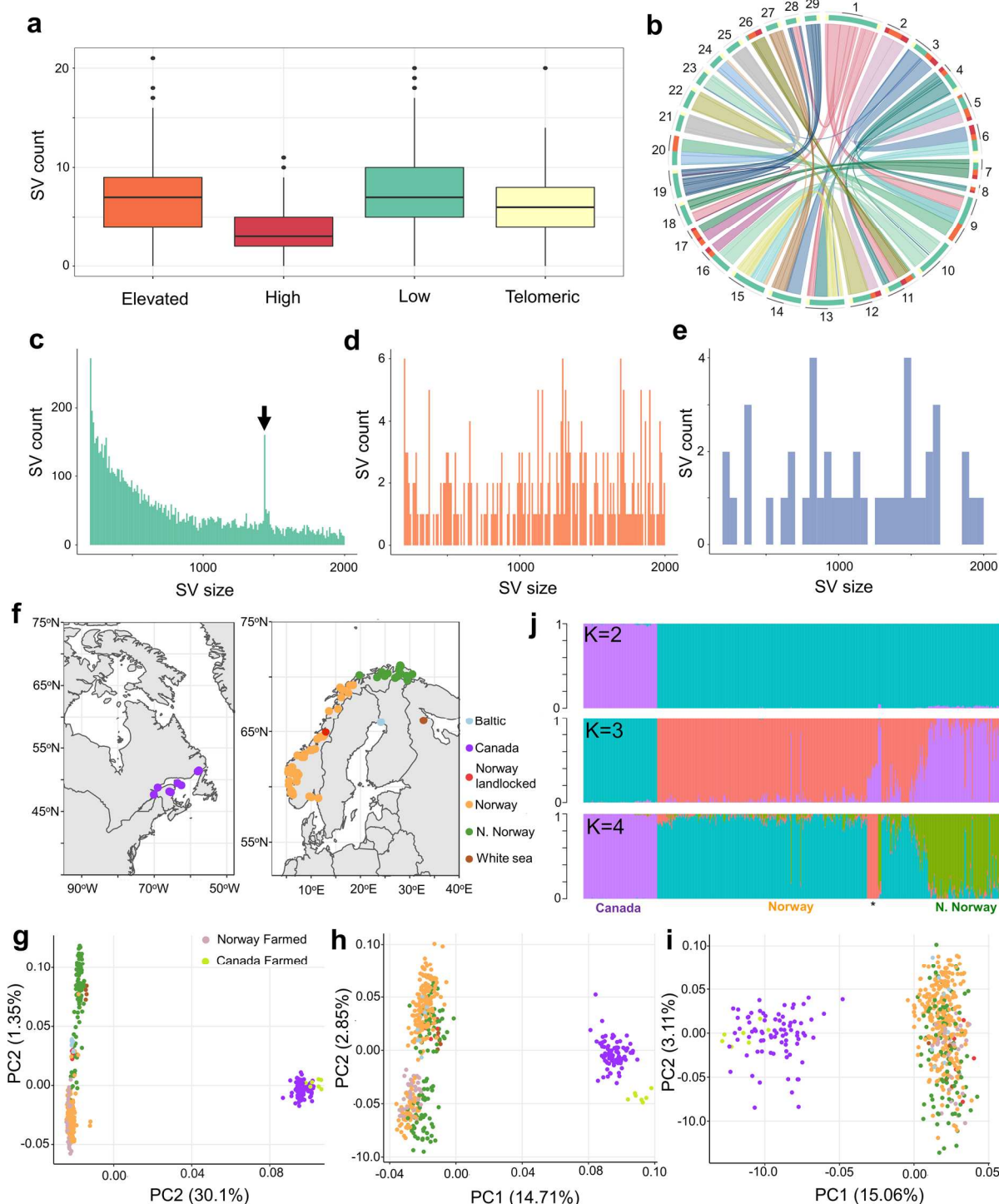
**Fig. 1**. SV landscape in 492 Atlantic salmon genomes. **a** SV counts per 1-million bp window in the genome split into homology categories[17] representing duplicated regions retained from the Ss4R WGD sharing 'low' (<90% identity), 'elevated' (90-95% identity) and 'high' (>95% identity) similarity in addition to telomere regions. **b** Locations of the same regions depicted on a Circos plot using the same colour scheme. **c-e** Size distributions of SVs for deletions (**c**), duplications (**d**) and inversions (**e**) with X axis limited to SVs ≤ 2,000 bp. Arrow in part **c** marks outlier peak in deletion calls (see Fig. 2). **f** Sampling locations of wild populations. **i-h** PCA of for each SV class: 14,017 deletions (**g**), 1,244 duplications (**h**), 242 inversions (**i**) with population matched by colour to part **f** for wild fish, and additional symbols given for farmed fish. **j** NGSadmix[86] analysis of 14,017 deletions with K=2, 3 and 4. Each individual is a vertical line with different colours marking genetically distinct groups. Asterisk corresponds to White sea, Baltic and landlocked populations (K=4 plot).
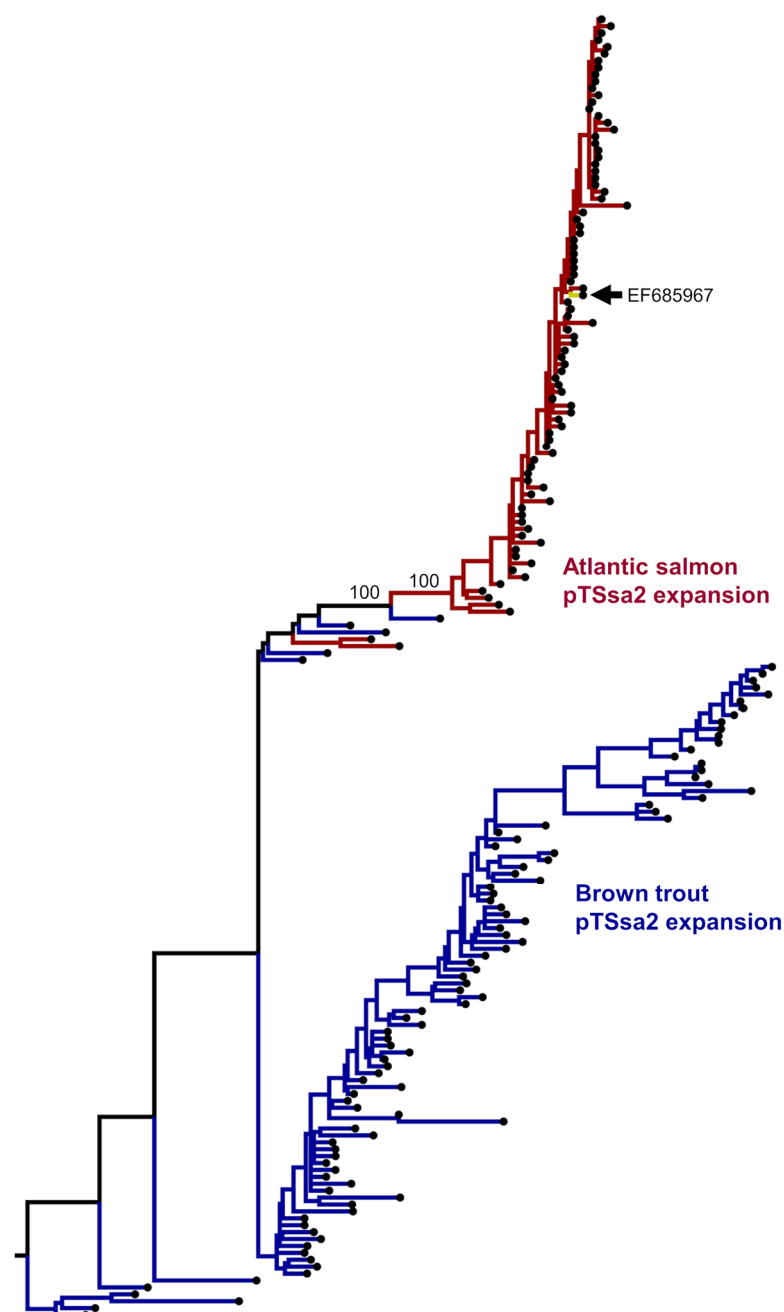
31

**Fig. 2**. Evidence for an active DNA transposon in *Salmo* evolution. Phylogenetic tree of Atlantic salmon sequences representing deletion polymorphisms matching the pTSsa2 piggyBac-like DNA transposon[28] (EF685967) and 100 top hits to this sequence within the brown trout genome. The tree was generated from an alignment spanning the length of pTSsa2 (Supplementary Data 3) using the TPM3+F+G4 substitution model. Bootstrap values are given at key nodes. A full tree with sequence identifiers, genomic locations of pTSsa2 sequences and bootstrap values is provided in Supplementary Fig. 18. A circos plot highlighting the location of pTSsa2 sequences in the Atlantic salmon genome is given in Supplementary Fig. 19.
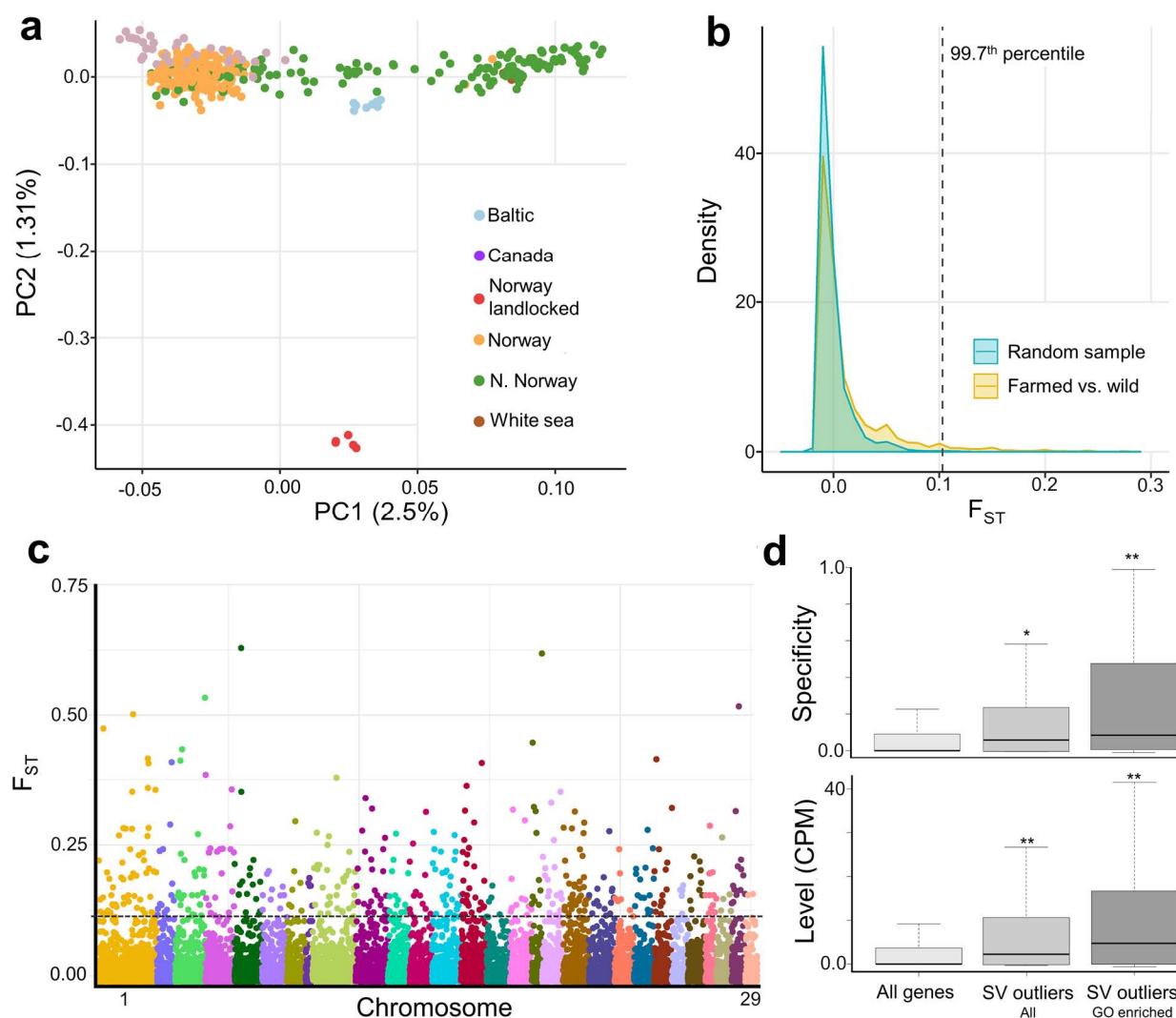
**Fig. 3.** Genetic differentiation of SVs between farmed and wild Atlantic salmon. **a** PCA used to select appropriate wild individuals for $F_{ST}$ comparison (n=257) vs. farmed salmon (n=34) on the basis of genetic distance by latitude (see also Supplementary Fig. 15) separated along PC1. The population symbols are the same as shown in Fig. 1. **b** Observed $F_{ST}$ value distribution comparing farmed vs. wild salmon contrasted against 200 random distributions for the same number of individuals. Dotted line shows cut-off $F_{ST}$ value employed in addition to a per SV criteria of $P<0.01$. **c** Manhattan plot of 12,627 $F_{ST}$ values with dotted line showing the same cut-off above which are the 584 SV outliers. **d** Brain gene expression specificity (top panel) and expression level (bottom panel) are increased for genes linked to the 584 outlier SVs, with the effect more pronounced for a 326 gene subset contributing to significantly enriched GO terms, compared to 44,469 genes in a multi-tissue transcriptome. Single and double asterisks indicate $P<0.005$ and $P<0.00001$, respectively. The observed increase in expression was specific to brain (plots for other tissues shown in Supplementary Fig. 22 and 23). Statistical analysis for all tissues shown in Supplementary Table 9.
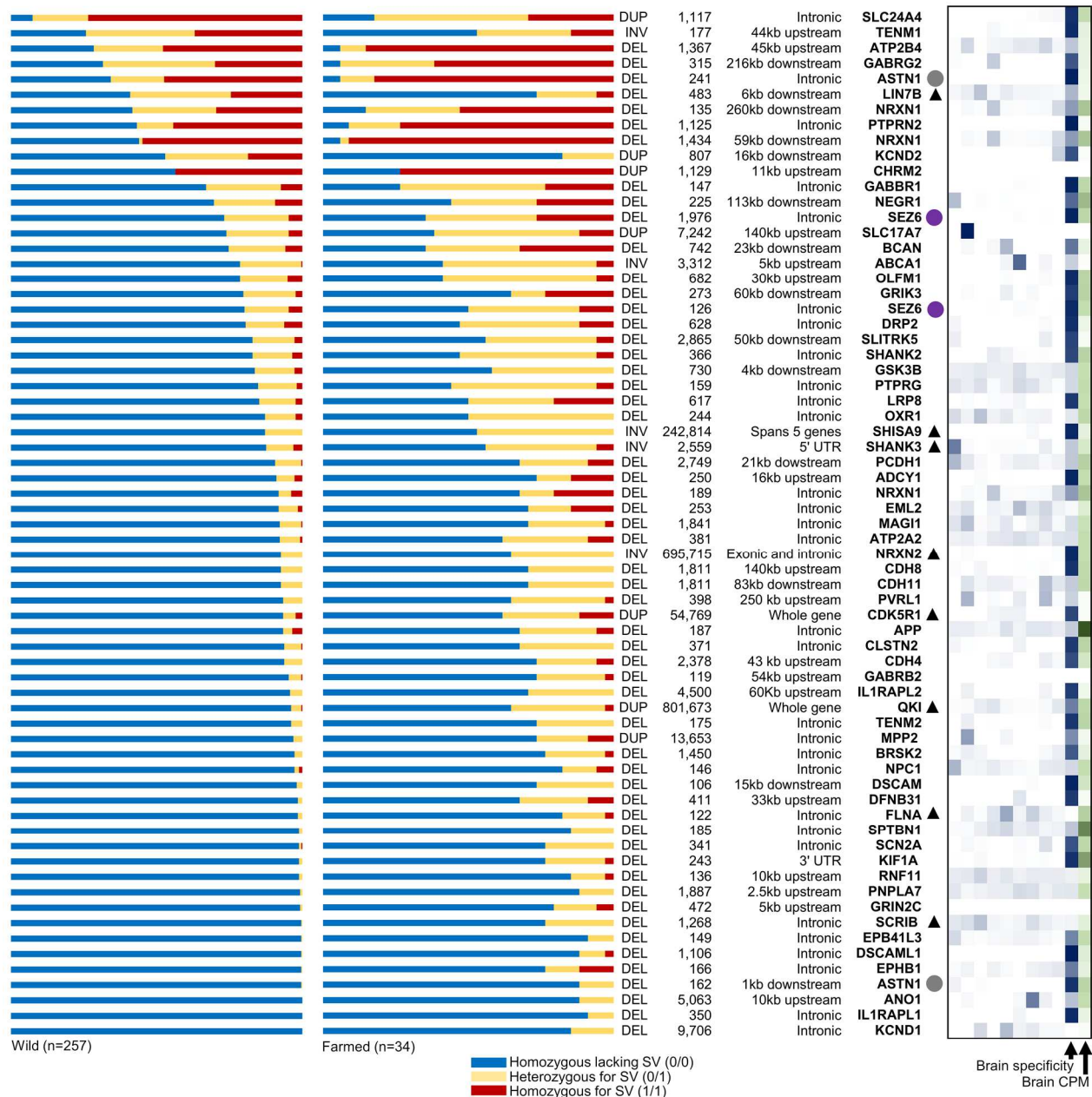
33

**Fig. 4**. SVs under selection during Atlantic salmon domestication are linked to 65 unique genes encoding synaptic proteins. SV genotypes are visualized on the left, ordered from bottom to top with decreasing frequency of homozygous genotypes (0/0) lacking the SV in wild fish. Annotation of each SV type, its size, and genomic location with respect to each synaptic gene is also shown. The circles next to genes highlight Ss4R ohnolog pairs and the black triangles indicate the overlap of an SV with a putative cis regulatory element (ATAC-Seq peak). The heatmap on the right depicts the expression specificity of each gene across an RNA-Seq tissue panel[17] (white to dark blue depicts lowest to highest tissue specificity; tissues shown in different columns from left to right: liver, gill, skeletal muscle, spleen, heart, foregut, pyloric caeca, pancreas, brain). The overall expression of each gene in brain is shown on the right of the heatmap (white to dark green depicts increasing CPM across the column). Data provided in Supplementary Table 10.

34

**Table 1**. Major effect SVs under divergent selection in farmed and wild Atlantic salmon

| | | | | | | SV genotype frequencies | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Chr** | **Start** | **Size** | **Type** | **Impact** | **$F_{ST}$** | **0/0 Wild** | **0/0 Farmed** | **0/1 Wild** | **0/1 Farmed** | **1/1 Wild** | **1/1 Farmed** |
| 1 | 15,177,232 | 23,362 | DEL | Deletes coding exons 3-12 in metabolic gene *SCCPDH* (LOC106569909, 12 exons) and lncRNA conserved in teleosts (LOC106569968) | 0.12 | 0.95 | 0.76 | 0.05 | 0.24 | 0.00 | 0.00 |
| 1 | 15,282,772 | 9,209 | DUP | Duplicates coding exons 5-10 within immune gene *GIMAP8* (LOC106569455, 14 exons) | 0.10 | 1.00 | 0.94 | 0.00 | 0.06 | 0.00 | 0.00 |
| 1 | 38,534,900 | 2,471 | DEL | Deletes coding exons 15-16 within sperm motility gene *CATSPERB* (106602505, 26 exons) | 0.11 | 1.00 | 0.91 | 0.00 | 0.09 | 0.00 | 0.00 |
| 1 | 53,229,610 | 801,673 | DUP | Duplicates region containing 9 coding genes, including immune gene *Pentraxin* (LOC100136583) | 0.27 | 0.96 | 0.65 | 0.04 | 0.32 | 0.00 | 0.03 |
| 1 | 63,072,912 | 1,133 | DEL | Deletes coding exons 16-17 within cell fusion gene *ADAM12* (LOC106607406, 23 exons) | 0.15 | 1.00 | 0.94 | 0.00 | 0.03 | 0.00 | 0.03 |
| 1 | 134,577,173 | 742 | DEL | Deletes lncRNA conserved in salmonids (LOC106567697) | 0.28 | 0.95 | 0.68 | 0.05 | 0.26 | 0.00 | 0.06 |
| 2 | 8,188,202 | 8,134 | DEL | Deletes coding exons 5-10 within glycoprotein gene *TUFT1* (LOC106575489, 16 exons) | 0.12 | 0.98 | 0.85 | 0.02 | 0.15 | 0.00 | 0.00 |
| 2 | 15,507,544 | 2,071 | DUP | Duplicates coding exons 12-15 within *HMCN1* (LOC106578676, 19 exons) | 0.24 | 0.28 | 0.00 | 0.16 | 0.00 | 0.56 | 1.00 |
| 2 | 45,905,818 | 49,351 | DEL | Deletes coding exons 1-25 of cellular adhesion gene *ITGAL* (106588084, 29 exons) | 0.11 | 0.95 | 0.76 | 0.05 | 0.24 | 0.00 | 0.00 |
| 2 | 51,645,286 | 1,172 | DEL | Deletion within coding exon 9 (frameshift) of endocytosis gene *SMAP1* (LOC100286439, 10 exons) | 0.15 | 1.00 | 0.91 | 0.00 | 0.09 | 0.00 | 0.00 |
| 3 | 53,262,801 | 56,833 | DUP | Disrupts coding sequence and intergenic region of two tandem *HEBP2* genes (LOC106600932, LOC106600932) | 0.19 | 0.96 | 0.79 | 0.04 | 0.12 | 0.00 | 0.09 |
| 4 | 33,772,841 | 2,115 | DEL | Deletes coding exons 21-26 of *PCNX1* (LOC106602984, 32 exons) | 0.24 | 1.00 | 0.91 | 0.00 | 0.03 | 0.00 | 0.06 |
| 5 | 23,514,943 | 157 | DEL | Deletes coding exon 8 of *PIGG* isoform 2 (LOC106604548, 8 exons) causing a frameshift | 0.35 | 1.00 | 0.76 | 0.00 | 0.24 | 0.00 | 0.00 |
| 5 | 29,459,708 | 1,886 | DEL | Deletes coding exons 2-3 within GTPase-activating gene *TBC1D2* (LOC106604634, 16 exons) | 0.10 | 1.00 | 0.94 | 0.00 | 0.06 | 0.00 | 0.00 |
| 5 | 54,982,436 | 5,313 | DUP | Affecting coding exons 6-8 within circadian regulator gene *NR1D2* (LOC100136378, 8 exons). Introduces stop codon | 0.15 | 0.84 | 0.50 | 0.10 | 0.38 | 0.06 | 0.12 |
| 6 | 1,542,320 | 19,710 | DUP | Duplicates coding exons 5-7 within immune gene *CD22* (106606237/8, 8 exons) | 0.13 | 0.87 | 0.62 | 0.10 | 0.29 | 0.03 | 0.09 |
| 6 | 29,579,766 | 5,320 | DEL | Deletes lncRNA conserved in salmonids (LOC106607070) | 0.20 | 0.85 | 0.53 | 0.14 | 0.35 | 0.01 | 0.12 |
| 7 | 21,191,252 | 422,735 | INV | Inverts region containing 16 coding genes | 0.11 | 1.00 | 0.91 | 0.00 | 0.09 | 0.00 | 0.00 |

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 9 | 21,282,095 | 11,299 | DUP | Duplicates coding exon 2 within *PGBD3* (LOC106611080, 4 exons) | 0.12 | 0.99 | 0.91 | 0.01 | 0.06 | 0.00 | 0.03 |
| 9 | 53,275,027 | 100,799 | DUP | Fusion of region containing last 10 coding exons of *TAPT1* (LOC106611550) with first 4 coding exons of *PROM1* (LOC106611549 | 0.15 | 0.84 | 0.56 | 0.12 | 0.29 | 0.03 | 0.15 |
| 10 | 23,225,394 | 32,774 | DEL | Deletes region containing six tRNA genes | 0.14 | 0.99 | 0.85 | 0.01 | 0.15 | 0.00 | 0.00 |
| 11 | 13,465,612 | 5,950 | DEL | Deletes exon 1 within lncRNA conserved in teleosts (LOC106562070, 3 exons) | 0.10 | 1.00 | 0.94 | 0.00 | 0.06 | 0.00 | 0.00 |
| 12 | 21,083,103 | 1,693 | DEL | Deletes coding exon 2-3 within uncharacterised gene (LOC106564648, 6 exons) | 0.25 | 0.96 | 0.71 | 0.04 | 0.24 | 0.00 | 0.06 |
| 14 | 14,287,987 | 18,976 | DUP | Duplicates coding exons 8-15 within melanosome transport gene *MYRIP* (LOC106568916, 15 exons) | 0.36 | 0.96 | 0.62 | 0.02 | 0.24 | 0.02 | 0.15 |
| 14 | 83,617,466 | 91,512 | DUP | Duplicates region containing 9 coding exons from *FAM126A* (LOC106570580), complete cytokine gene *IL6* (LOC106570581) and coding exon 1 from *RAPGEF5* (LOC106570584) | 0.13 | 0.98 | 0.88 | 0.02 | 0.06 | 0.00 | 0.06 |
| 18 | 56,889,482 | 39,099 | DUP | Duplicates coding exons 1-12 within immune gene *CD22* (LOC106577812, 20 exons) | 0.12 | 0.94 | 0.76 | 0.05 | 0.18 | 0.01 | 0.06 |
| 18 | 64,338,324 | 852 | DEL | Deletes coding exon 7 within gene *PARP14*-like (LOC106578007, 7 exons) and ablates stop codon | 0.15 | 0.84 | 0.56 | 0.14 | 0.32 | 0.02 | 0.12 |
| 19 | 51,422,161 | 31,121 | INV | Flips coding exon 1-2 within fatty acid elongation gene *ELOVL6* (LOC106579283, 4 exons) | 0.11 | 0.93 | 0.71 | 0.07 | 0.29 | 0.00 | 0.00 |
| 22 | 40,200,901 | 5,863 | DEL | Deletes coding exon 2 within *PLEKHA6* (LOC106583501, 24 exons) | 0.13 | 0.97 | 0.85 | 0.02 | 0.06 | 0.01 | 0.09 |
| 24 | 11,833,364 | 165 | DUP | Deletes half of coding exon 2 within tRNA methyltransferase gene *TRMT2A* (LOC106584929, 12 exons) | 0.11 | 0.09 | 0.32 | 0.44 | 0.47 | 0.46 | 0.21 |
| 24 | 19,661,320 | 266,147 | INV | Affects 6 coding genes, inverting 5 genes completely and all but first exon of *AAK1* (LOC106585601) | 0.16 | 0.83 | 0.44 | 0.17 | 0.56 | 0.00 | 0.00 |
| 27 | 42,220,948 | 341 | DEL | Partially deletes exon 4 in angiogenesis gene (*ANG2* LOC106589146, 5 exons) causing frameshift | 0.12 | 0.56 | 0.24 | 0.34 | 0.50 | 0.10 | 0.26 |
| 28 | 3,887,040 | 5,373 | DUP | Duplication affecting zinc transporter gene *SLC39A11* (LOC100380452, 10 exons) causing frameshift | 0.14 | 0.95 | 0.79 | 0.04 | 0.09 | 0.02 | 0.12 |
| 28 | 16,046,880 | 24,780 | DUP | Fusion involving coding exons 9-16 of sodium transport gene *SLC38A10* (LOC106589592, 16 exons) and exons 1-3 of vesicular transport gene *TEPSIN* (15 exons) | 0.52 | 0.86 | 0.29 | 0.10 | 0.26 | 0.04 | 0.44 |

Genotypes: 0/0: homozygous lacking SV; 0/1 heterozygous for SV 1/1 homozygous for SV