1    **Dissecting the loci underlying maturation timing in Atlantic salmon using haplotype and multi-SNP**

2    **based association methods**

3

4    Marion Sinclair-Waters[1,2], Torfinn Nome[3], Jing Wang[3,4], Sigbjørn Lien[3], Matthew P. Kent[3], Harald Sægrov[5],

5    Bjørn Florø-Larsen[6], Geir H. Bolstad[7], Craig R. Primmer[1,2*] & Nicola J. Barson[3*]

6

7    [1]Organismal and Evolutionary Biology Research Programme, Faculty of Biological and Environmental

8    Sciences University of Helsinki, Helsinki, Finland

9    [2]Institute of Biotechnology, Helsinki Institute of Life Sciences, University of Helsinki, Helsinki, Finland

10   [3]Centre for Integrative Genetics, Department of Animal and Aquacultural Sciences, Faculty of Biosciences,

11   Norwegian University of Life Sciences, Ås, Norway

12   [4]Key laboratory for Bio-Resources and Eco-Environment, College of Life Science, Sichuan University,

13   Chengdu, China

14   [5]Rådgivende Biologer, Bergen, Norway

15   [6]Norwegian Veterinary Institute, Trondheim, Norway

16   [7]Norwegian Institute for Nature Research (NINA), Trondheim, Norway

17   *Shared last authors

18

19 **ABSTRACT**

20 Resolving the genetic architecture of fitness-related traits is key to understanding the evolution and

21 maintenance of fitness variation. However, well-characterized genetic architectures of such traits in wild

22 populations remain uncommon. In this study, we used haplotype-based and multi-SNP Bayesian association

23 methods with sequencing data for 313 individuals from wild populations to further characterize known

24 candidate regions for sea age at maturation in Atlantic salmon (*Salmo salar*). We detected an association at

25 five loci (on chromosomes *ssa06*, *ssa09*, *ssa21*, and *ssa25*) out of 116 candidates previously identified in an

26 aquaculture strain with maturation timing in wild Atlantic salmon. We found that at each of these five loci,

27 variation explained by the locus was predominantly driven by a single SNP suggesting the genetic

28 architecture of Atlantic salmon maturation includes multiple loci with simple, non-clustered alleles. This

29 highlights the diversity of genetic architectures that can exist for fitness-related traits. Furthermore, this study

30 provides a useful multi-SNP framework for future work using sequencing data to characterize genetic

31 variation underlying phenotypes in wild populations.

32 INTRODUCTION

33     Understanding the genetic processes underlying fitness variation is a fundamental goal in evolutionary

34 biology. Identifying genetic variants that underlie fitness-related traits is therefore crucial, yet remains

35 challenging. Substantial effort has been made to characterize the genetic architecture of traits – i.e. Are there

36 few or many loci involved? Are loci effects small or large? How are loci distributed across the genome? And

37 what are the allele frequencies at these loci [1–5]? It is generally assumed that in most cases single genetic

38 variants translate into only small changes in complex traits, and therefore follow a polygenic [6,7] or an

39 omnigenic [3,8] model of inheritance.

40     Among genome-wide association studies published to date, many complex traits appear to be

41 polygenic [9]. Although polygenicity is widespread, an increasing number of examples of major effect loci

42 exist, whereby one locus explains a large proportion of the phenotypic variation [10,11]. In some cases,

43 major effect loci can contain multiple tightly linked genes, coined "supergenes", where localized reduction in

44 recombination is often caused by larger chromosomal rearrangements. For example, this phenomenon is

45 known to underlie phenotypic variation observed among ruff (*Philomachus pugnax*) mating morphs [12,13],

46 Atlantic cod (*Gadus morhua*) [14,15] and rainbow trout (*Oncorhynchus mykiss*) migratory ecotypes [16], and

47 *Heliconius* butterfly wing-pattern morphs [17]. More recent work has found that major effect loci can exist

48 alongside a polygenic background where loci with a variety of effect sizes underlie trait variation [18,19].

49 Such mixed genetic architectures may be pervasive, but currently remain undetected due to the large sample

50 sizes required for detecting loci with smaller effects [19] and it is possible that additional examples are to be

51 found with future higher-powered studies. Although studies aimed at resolving genotype-phenotype links are

52 mounting, well-characterized genetic architectures of fitness-related traits, particularly in natural populations,

53 are still uncommon.

54     While some trait-associated loci have been identified, such findings lead to other crucial questions:

55 How have trait-locus associations arisen? Has the locus arisen through a single or multiple new mutations?

56 Or alternatively, did the locus emerge via recombination that gave rise to new combinations of existing

57 variants? Numerous studies from the past decade have shown that major effect loci involve the cumulative

58 effects of multiple mutations, rather than a single mutation, thus highlighting the relevance of considering the

59    latter scenarios. For example, Bickle et al. [20] found that ~60% of variation in female abdominal

60    pigmentation in *Drosophila melanogaster* can be explained by sequence variation at the *bab* locus, but a

61    GWAS (genome-wide association study) analyzing the same trait did not identify a single SNP in *bab* that

62    passed the genome-wide significance threshold. Alleles consisting of multiple SNPs were associated with

63    high proportions of the variation, whereas, single SNPs had only small effects and were therefore missed in

64    the single-SNP GWAS. Additionally, Linnen at al. [11] and Kerdaffrec et al. [21] also identify multiple

65    mutations within a confined region that have cumulative effects on colour traits in deer mice and seed

66    dormancy in *Arabidopsis thaliana*, respectively. In natural populations with gene flow such as in Linnen et

67    al. [11] and Kerdaffrec et al. [21], this is perhaps not unexpected as theory predicts that clustered and major

68    effect loci will evolve under such scenarios [22,23]. Given these findings, examining extended sequence

69    haplotypes containing multiple SNPs, rather than each SNP independently, is important [24]. This can be

70    achieved by using alternative strategies that look at combined effects of variants, rather than single-SNP

71    methods typically used in GWAS.

72        Here we investigate the genetic basis of Atlantic salmon (*Salmo salar*) sea age at maturity – the

73    number of years spent in the marine environment before reaching maturity and returning to the natal river

74    (freshwater) to reproduce. Age at maturity is an important life history trait affecting fitness traits such as

75    survival, size at maturity and reproductive success [25,26]. Substantial variation in Atlantic salmon sea age

76    at maturity is maintained due to a trade-off between mating success at spawning grounds and survival,

77    whereby individuals that mature later are larger and have higher reproductive success on the spawning

78    grounds, but lower survival and thus lower chance of reaching reproductive age. In contrast individuals that

79    mature early are smaller and have lower reproductive success, but higher survival and thus higher chance of

80    reaching reproductive age [27,28].

81        Variation in maturation timing in Atlantic salmon is highly heritable [19,29,30] and consequently

82    there is substantial interest in understanding the underlying genetic architecture. A large-effect locus on

83    chromosome 25 explaining up to 39% of the variation in sea age at maturity was found in wild European

84    populations [10] and domesticated salmon [31]. The primary candidate gene underlying the association of

85    this locus is *vgll3* due to its close proximity to the associated SNP variation [10,31,32] and its known

4

86    function in other species. The *vgll3* gene encodes a transcription cofactor that, amongst other things,

87    regulates adipogenesis [33] and is associated with variation in puberty timing in humans [34,35]. In addition

88    to *vgll3*, Sinclair-Waters et al. [19] identified 119 other candidate genes for male maturation in a GWAS

89    including >11,000 males from the same Atlantic salmon aquaculture strain. Two particularly strong

90    associations between maturation timing were found on chromosome 9 in close proximity to *six6* and

91    chromosome 25, *vgll3*. The association of *six6* was also found by Barson et al. [10] in wild Atlantic salmon,

92    however, the signal disappeared after correction for population structure. Interestingly, the *six6* gene is also

93    associated with age at maturity in two Pacific salmon species [36], humans [35] and cattle [37]. However,

94    Barson et al. [10] focused solely on single-SNP associations via GWAS without considering the possible

95    influence of combined variant effects.

96        Studies using sequencing data to examine variation associated with important fitness-related traits in

97    wild populations are limited. However due to developments in sequencing technologies and bioinformatics,

98    studies using this approach are likely to rise in number. We therefore aim to provide a useful and timely

99    framework for characterizing genetic variation underlying phenotypes in wild populations in the future.

100   Here, we focus on further characterizing the association between the loci identified in Sinclair-Waters et al.

101   [15] and sea age at maturity in wild Atlantic salmon. We integrate re-sequencing data and phenotype

102   information for 313 individuals from 53 wild population of Atlantic salmon with alternative GWAS

103   strategies that consider the combined effects of variants, rather than single-SNP effects. This approach can

104   provide better resolution of the variants that are potentially involved in controlling fitness-related traits such

105   as maturation timing in Atlantic salmon.

106

107   METHODS

108   *Study material*

109       Whole genome sequencing data was obtained for 313 wild individuals collected from 53 Norwegian

110   and Finnish populations spanning the Norwegian coast and to the Barents sea in the north (59°N - 71°N)

111   (Supplementary Table S1) previously reported in Bertolotti et al. [38]. The 313-individual dataset includes

112 populations belonging to both the Atlantic and Barents/White sea phylogeographic groups. These regions

113 were studied in Barson et al. [10] using SNP-array data and a single SNP approach, therefore missing

114 variants and potentially combined variant effects. Individuals were categorized into three maturation

115 categories based on the number of years spent at sea prior to their first return migration to rivers for

116 spawning: 1 (one year spent at sea), 2 (two years spent at sea), or 3 (three or more years spent at sea). Only

117 five individuals had spent four years and were therefore combined with three-year fish for all analyses.

118 *SNP calling & filtering*

119 Variant calling and the first round of filtering was done in a larger set of individuals described in

120 Bertolotti et al. [38]. Raw Illumina reads were mapped to the Atlantic salmon genome (ICSASG_v2) [39]

121 using *bcbio-nextgen v.1.1* [40]with the *bwa-mem aligner v.0.7.17* [41]. Genomic variation was identified

122 using the Genome Analysis Toolkit (*GATK*) *v4.0.3.0.*, following *GATK*'s best practice recommendations.

123 *Picard v2.18.7* [42] was used to mark duplicates and *GATK* was used for joint calling [43]. Variants were

124 annotated using *SNPeff v. 4.3* [44]. Variant call were further filtered with GATK's variant filtration

125 according to the following *--filterExpression*: "MQRankSum < -12.5 || ReadPosRankSum < -8.0 || QD < 2.0

126 || FS > 60.0 || (QD < 10.0 && AD[0:1] / (AD[0:1] + AD[0:0]) < 0.25 && ReadPosRankSum < 0.0) || MQ <

127 30.0". SNPs were then filtered using *SNPable* procedure [45], where 100 bp kmers are mapped to reference

128 genome (ICSASG_v2) using Burrows-Wheeler Aligner (*bwa aln*) [46], and only SNPs within regions with

129 reads that uniquely map are retained. We then removed additional SNPs with *vcftools* using the following

130 criteria: *--min-alleles 2, --max-alleles 2, --maf 0.0000000001, --max-missing 0.7, --remove-indels, --minGQ*

131 *10,* and *–minDP 4*. A subset 313 individuals from wild populations was then extracted from this larger

132 dataset using *vcftools* [47]. This reduced dataset was used for all subsequent analyses.

133 *Principal component analysis*

134 We produced a reduced SNP dataset by pruning one SNP from each SNP pair with a correlation

135 coefficient ($r^2$) greater than 0.2 within a 50 kb block using the *--indep-pairwise 50 10 0.2* function

136 implemented in *PLINK v1.9* [48]. This yielded 403,540 SNPs to examine population structure using a

137 principal component analysis, *smartpca,* implemented in the EIGENSOFT *v5* software [49].

138    *Data preparation*

139        In this study, we focus on genomic regions containing the 116 candidate loci for age at maturity

140    identified in Sinclair-Waters et al. [19]. We extracted SNP genotype data from 500 kb regions surrounding

141    the 116 trait-associated SNPs identified in Sinclair-Waters et al. [19] using *vcftools'* [47] position filtering

142    functions *--from-bp* and *--to-bp,* as well as allele filtering function *--mac 1* to keep only polymorphic sites.

143    SNPs that were within 250 kb of an adjacent SNP were analyzed together by examining a region that extends

144    250 kb upstream of the first SNP to 250 kb downstream of the last SNP.

145        The current Atlantic salmon genome (ICSASG_v2) contains a known assembly error within the 500

146    kb region surrounding the known candidate loci *vgll3* [31]. A misplaced and misoriented scaffold currently

147    placed downstream of *vgll3* belongs within a gap in the assembly just upstream of *vgll3* on ssa25. For this

148    reason, we constructed a revised assembly for this chromosome. SNP calling was performed as described

149    above. We then retained SNPs that had met the filtering criteria. A total of 8 candidate SNPs are located

150    within regions of the genome that were moved. To find the position of these SNPs in the revised

151    chromosome 25 sequence, we extracted 200 bp surrounding each of these SNPs from the current genome

152    assembly (ICSASG_v2) using the *getfasta* function in *BEDTools* [50]. The 200 bp sequence was then blasted

153    to the fixed assembly to determine the new position of each SNP using Blast's *blastn* function [51]. Using

154    the new SNP positions, SNP genotypes within a 500 kb region surrounding the moved candidate SNPs were

155    extracted from the fixed dataset using *vcftools*.

156    *Association testing at candidate regions*

157        We applied three association mapping methods to describe the genetic architecture underlying sea age

158    at maturity at each of the candidate regions identified in Sinclair-Waters et al. [19]. First, a multi-SNP

159    approach examining associations between phenotype and haplotypes was conducted using Bayesian linear

160    regression implemented in *hapQTLv1.00* [52]. In this approach, a hidden Markov model is used to

161    characterize haplotype structure and ancestry [53]. Haplotype sharing at each marker is then used to quantify

162    genetic similarity among individuals. Haplotype associations are identified by testing for an association

163    between genetic similarity at each marker and the phenotype [52]. Each of the extracted *vcf* files was

7

164     converted to *bimbam* format using *PLINK 1.9* [54]. The resulting *bimbam* files were used as input for

165     *hapQTL.* Second, single SNP associations were also identified using a Bayesian linear regression method

166     implemented in *hapQTL* [55]. For all *hapQTL* association tests, sex and the six most significant principal

167     components (see above) were included as covariates in the models. Each *hapQTL* run consisted of 2 EM runs

168     (-e 2) with 40 steps (-w 40), 2 upper clusters (-C 2), 10 lower clusters (-c 10).  Three replicate *hapQTL* runs

169     were performed for each of the 116 selected regions. Based on recommendations from Jeffreys [56], Bayes

170     factors greater than three were considered evidence for an association of either SNPs or haplotype with sea

171     age at maturity phenotype.

172         Third, a multi-SNP approach aimed to estimate the number and identity of SNPs underlying trait

173     variation at each candidate region using Bayesian Variable Selection regression implemented in *PiMASS*

174     [55]. Due to computational restrictions, the *PiMASS* analysis was performed for only candidate regions that

175     had a SNP or haplotype association with Bayes factor greater than 3. Prior to the *PiMASS* analysis, all

176     missing genotypes were imputed in BIMBAM [55] as mean genotypes (-wmg) using default settings.

177     Additionally, our phenotype values for sea age at maturity were adjusted to correct for confounding effects

178     of sex and population structure by regressing the phenotype on sex and the six most significant principal

179     components (see above) using the *lm* function in *R*. *PiMASS* was run with the residual phenotype values. We

180     placed priors on the proportion of variance explained by SNP(s) (hmin = 0.001 and hmax = 0.999) and the

181     number of SNPs in the model (pmin = $\log\frac{1}{N}$ and pmax = $\log\frac{300}{N}$, where N is the total number of SNPs). Each

182     run consisted of a burn-in of 1000000 steps, followed by 2500000 steps where parameter values were

183     recorded every 1000 steps. For each analysis, we examined the posterior inclusion probability for each SNP,

184     the distribution of the number of included SNPs and the distribution of the proportions of variance explained

185     per model. We also examined the path of estimated Bayes factors and parameter values (h, p, s) across all

186     recorded iterations to check for convergence of runs.

187         To further assess whether more than one SNP in a candidate region was significantly associated with

188     sea age at maturity, we regressed out the top-associated SNP from the residual phenotype values described

189     above and reran *PiMASS* using the previously-used priors and settings. We then examined the posterior

190     inclusion probability for each SNP, the distribution of the number of included SNPs, and the distribution of

8

191    proportion of variance explained to determine whether there was evidence for multiple SNP associations

192    within a given candidate region.
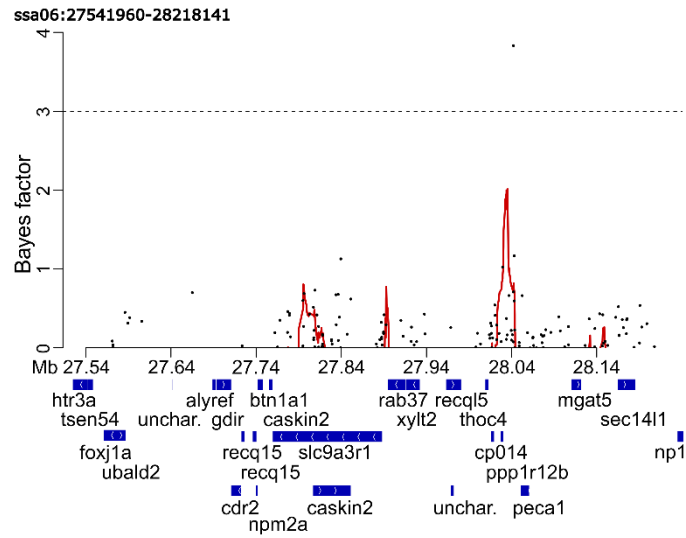
193

194    RESULTS

195    *Principal component analysis*

196        The first six principal components (PCs) calculated with the pruned SNP dataset explained 1.96%,

197    0.68%, 0.63%, 0.59%, 0.56% and 0.51% of the genetic variance, respectively (Supplementary Figure S1).

198    These six PCs were included in subsequent association analyses to reflect population structure among

199    samples.

200    *Associations identified with hapQTL*

201        Single-SNP and haplotype association analyses with *hapQTL* revealed strong (Bayes factor > 3)

202    association signals at 5 of the 116 candidate regions (Figure 1, Supplementary Figure S2). The strongest

203    association observed within each region was with a single SNP, rather than an extended haplotype,

204    suggesting a single mutation underlies the effect of each of these regions on maturation timing. However,

205    exceptions occurred in the ssa09:24636574-25136574 and ssa25:28389273-28889273 regions, where second

206    association signals were found upstream of the primary association signal and were most strongly linked to

207    an extended haplotype. For instance, strong haplotype association scores (Bayes factor > 3) spanned a 26971

208    bp region (ssa09:24781742-24808713) containing an uncharacterized gene (LOC106610978) and *pcnx4*. In

209    the ssa25:28389273-28889273 region, a strong haplotype signal was found within *edar* (Figure 1).
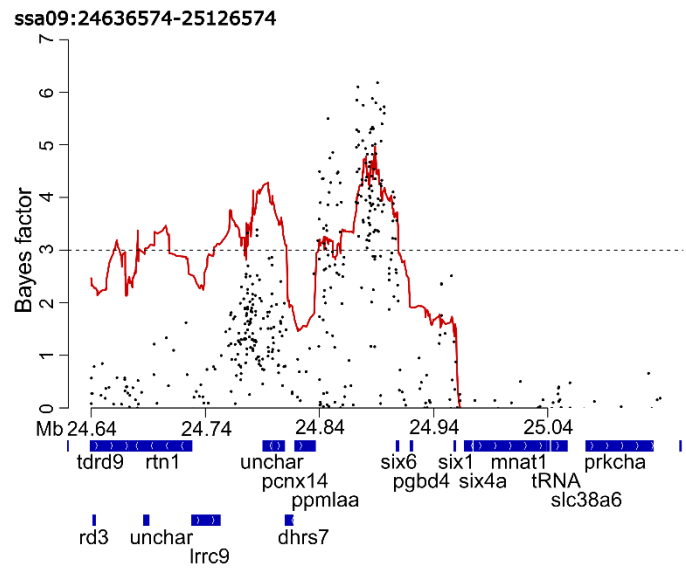
210        We find differences in the location of the top-associated SNPs found here and those identified in

211    Sinclair-Waters et al. [19]. For regions ssa06:27541960-28218141, ssa09:10915066-11415066 and

212    ssa25:28389273-28889273, the top-associated SNP was located further upstream than in Sinclair-Waters et

213    al. [19]. Contrastingly, the strongest associated SNPs within the regions ssa09:24636574-25136574 and

214    ssa21:49390687-49890687 differed only slightly (<5000 bp) between studies (Table 1).
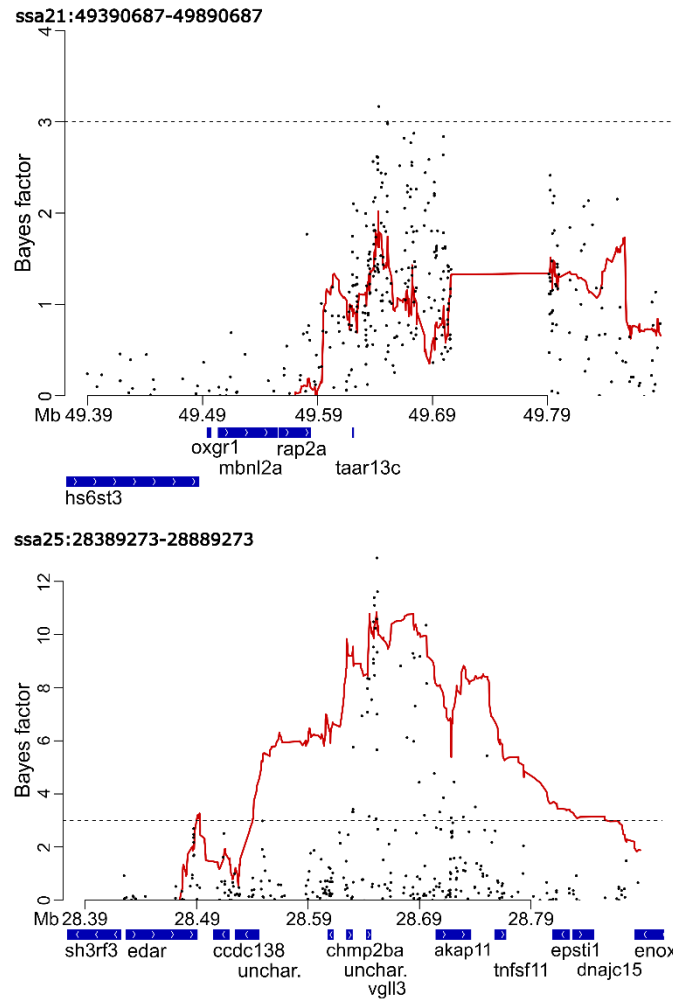
ssa06:27541960-28218141



ssa09:10915066-11415066

215

216



ssa09:24636574-25126574

217

10

Figure 1. Plots displaying single SNP associations (black points) and haplotype associations (red line) scores from *hapQTL* for the five candidate regions with Bayes factors greater than 3. Y-axis shows the Bayes factor indicating the association strength. X-axis shows the position on the respective chromosomes.

223    Table 1. Strongest association signals for each candidate region showing evidence of an association with sea age at maturity, the genes in closest proximity
224    and association values from *hapQTL*. Top SNPs for each region from previous SNP-array study [19].
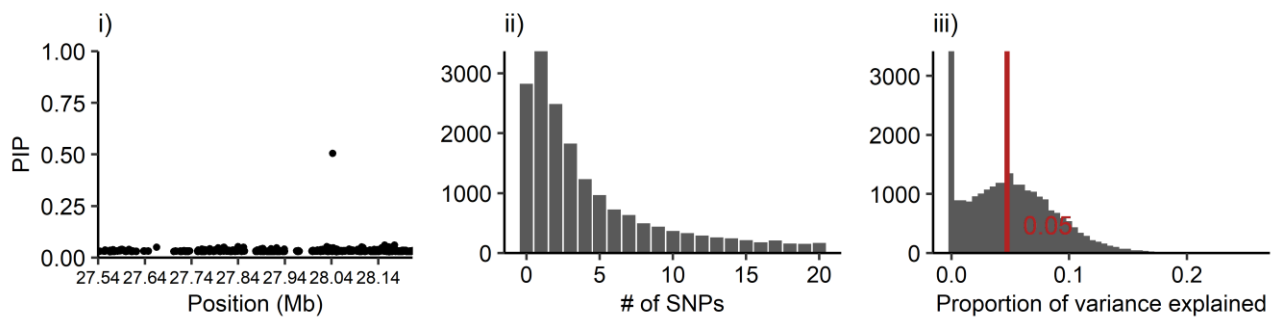
| Candidate region | Top signal | Closest gene | Bayes Factor | -log$_{10}$(*P*-value) | Allele frequency | Top SNP(s)[a] | Candidate gene(s)[a] |
|---|---|---|---|---|---|---|---|
| ssa06:27541960-28218141 | 6:28045390 (SNP) | *pecam1* (intron) | 3.835 | 5.107 | 0.320 | 6:27791960 6:27968141 | *slc9a3r1* *recql5* LOC106606978 |
| ssa09:10915066-11415066 | 9:11266848 (SNP) | *asap2a* (upstream) | 4.696 | 5.434 | 0.074 | 9:11165066 | *mboat2* |
| ssa09:24636574-25136574 | 9:24888841 (SNP) | *six6* (upstream) | 6.184 | 4.242 | 0.425 | 9:24886574 | *six6* |
| ssa21:49390687-49890687 | 21:49645222 (SNP) | *taar13c* (upstream) | 3.172 | 4.649 | 0.464 | 21:49640687 | *taar13c* |
| ssa25:28389273-28889273 | 25: 28651640 (SNP) [ICSASG_v2: 25:28669350] | *vgll3* (downstream) | 12.893 | 6.406 | 0.358 | 25:28910202 | *vgll3* |

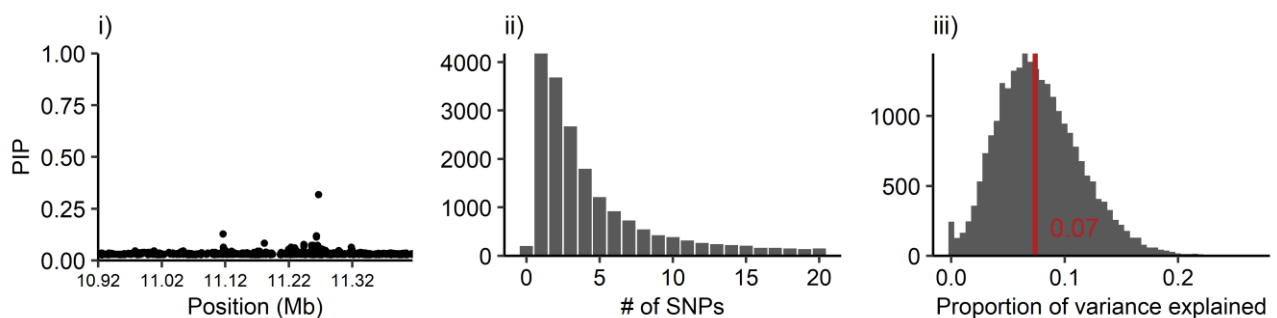225    [a]From Sinclair-Waters et al. [19].

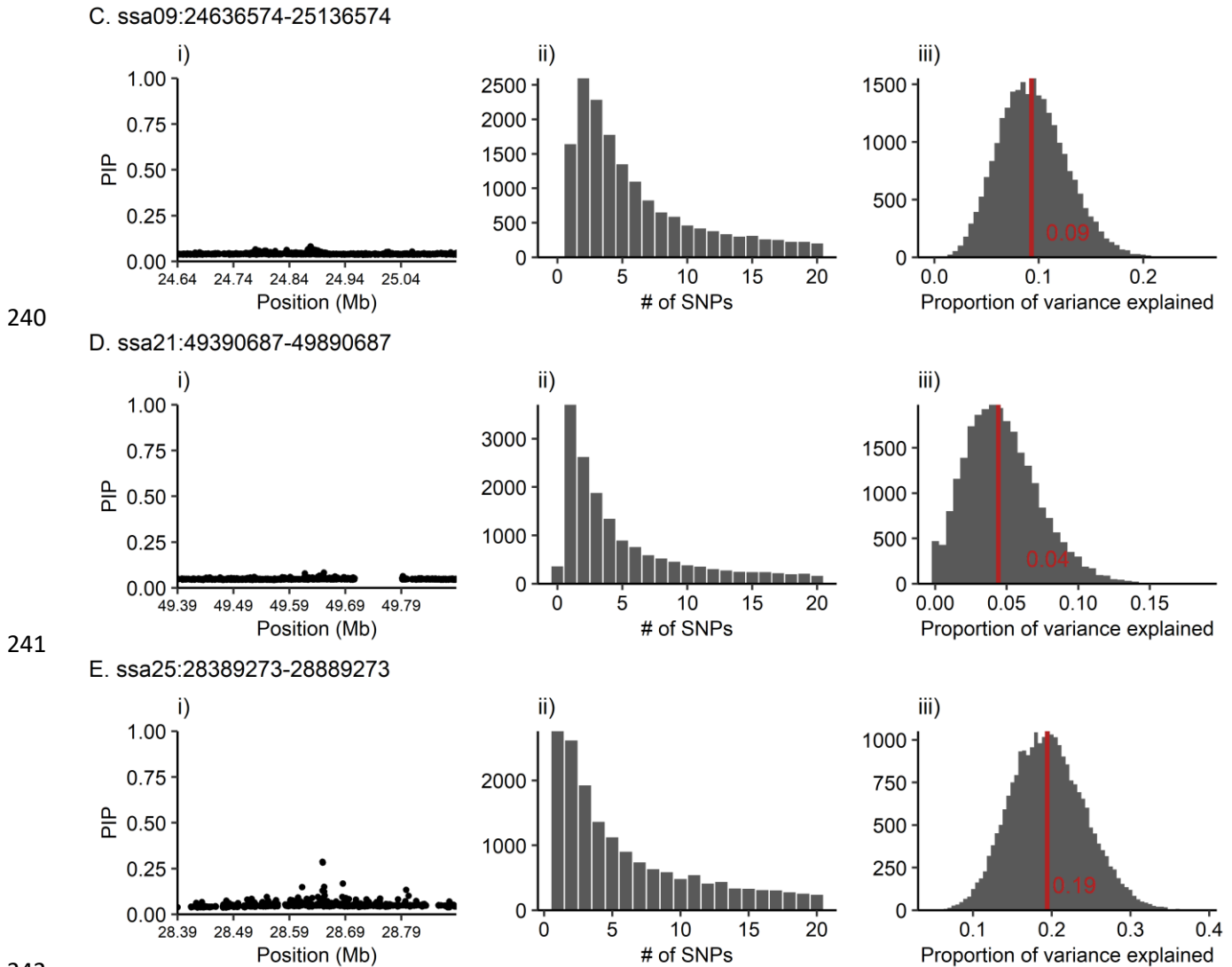226    *Multi-SNP associations identified using PiMASS*

227    Multi-SNP association analysis with *PiMASS* showed that at four of five candidate regions, a single-

228    SNP model was most commonly used to explain variation in sea age at maturity. At one candidate region,

229    ssa09:24636574-25136574, a multi-SNP model including two SNPs was most commonly used to explain

230    variation in sea age at maturity. Median proportion of variance explained by each candidate region ranged

231    between 4% and 19% (Figure 2, Table 2). However, when the top-associated SNP was regressed out from

232    the phenotype values, no SNPs were selected to explain sea age at maturity for all five candidate regions.

233    Additionally, post-regression median proportion of variance was substantially lower – ranging between 0%

234    and 1% (Supplementary Figure S3, Table 2). This would suggest that sea age variation explained by each of

235    these regions is largely driven by a single mutation. We observe no obvious trends in parameter values or

236    Bayes factors, suggesting models converged and burn-in period was adequate (Supplementary Figure S4

237    &S5).

A. ssa06:27541960-28218141



238

B. ssa09:10915066-11415066



239

C. ssa09:24636574-25136574



D. ssa21:49390687-49890687



E. ssa25:28389273-28889273



Figure 2. *PiMASS* results for each of the tested candidate regions: A. ssa06:27541960-28218141, B. ssa09:10915066-11415066 C. ssa09:24636574-25136574, D. ssa21:49390687-49890687, and E. ssa25:28389273-28889273. Plots display the following results for each candidate region: i) posterior inclusion probability (PIP) indicating the probability of a SNP being included in a model explaining sea age at maturity variation, ii) truncated distribution of the number of SNPs included in a model explaining sea age at maturity variation, and iii) distribution of proportion of variance explained per recorded iteration (2500). Red line indicates the median proportion of variance explained.

251 Table 2. *PiMASS* results prior to and after regression of top-associated SNP identified in the initial *PiMASS*
252 analysis. These include the mode of the distribution of the number of SNPs and the median of the
253 distribution of proportion of variance explained (PVE) for a model explaining sea age at maturity.

| Candidate region | Mode # of SNPs | Median PVE | Mode # of SNPs (post-regression) | Median PVE (post-regression) |
|---|---|---|---|---|
| ssa06:27541960-28218141 | 1 | 0.05 | 0 | 0 |
| ssa09:10915066-11415066 | 1 | 0.07 | 0 | 0.01 |
| ssa09:24636574-25136574 | 2 | 0.09 | 0 | 0.01 |
| ssa21:49390687-49890687 | 1 | 0.04 | 0 | 0 |
| ssa25:28389273-28889273 | 1 | 0.19 | 0 | 0.01 |

254

255

256 DISCUSSION

257     Despite that combined effects of multiple variants at trait-associated loci are playing an important role

258 in controlling fitness traits across a variety of species [11,20,21], our results indicate that sea age at

259 maturation in Atlantic salmon is predominantly associated with single SNP variation at candidate regions.

260 Using resequencing data to analyse 116 candidate loci and an analytical framework aimed at detecting multi-

261 SNP associations, we find that single SNPs explain the variation in sea age at maturity in almost all cases.

262 This work targeting candidate genes identified in aquaculture salmon strains suggests a mixed genetic

263 architecture where a combination large-effect loci and smaller-effect loci also underlies age at maturity in

264 wild Atlantic salmon populations. Two core loci, *vgll3* and *six6*, likely play a key role in determining age at

265 maturity and additional smaller effect loci may be important for fine-tuning the trait across heterogeneous

266 environments.

267     Theoretical modelling predicts that clustering of tightly linked adaptive mutations will occur under

268 gene flow and selection in populations inhabiting spatially and/or temporally heterogeneous environments

269 [22,23]. Although this seems to be a plausible scenario under which the genetic architecture of age at

270 maturity has evolved in Atlantic salmon, our work suggests that the association in each of the candidate

271 regions is driven by a single mutation. We cannot rule out, however, the possibility that the examined

272 regions have pleiotropic effects and contain SNPs controlling other adaptive traits that have weak or no

273 correlation with maturation timing. It is also possible that we did not have sufficient power to detect

15

274    additional SNPs in these regions with small effects or with rare alleles. However, previous empirical studies

275    have found few, but complex, loci with clusters of adaptive mutations [11,20,21], thus motivating our

276    investigation of multi-SNP and haplotypic effects. Remington [24] also highlights the importance of

277    distinguishing between allelic effects and single mutational effects when examining the genetic architecture

278    of adaptive variation and its evolution. Our findings, however, suggest that alternative genetic architectures

279    are feasible. One possible explanation could relate to the multiple whole genome duplication events that have

280    occurred in Atlantic salmon and other salmonids [57]. The presence of multiple gene copies may impact the

281    evolution of genetic architecture for traits such as age at maturity in Atlantic salmon. It is also possible that

282    gene flow among Atlantic salmon populations is too restricted to neighbouring populations and/or strength of

283    selection is insufficient for the establishment of linked mutations, as there is a rather specific balance of gene

284    flow and selection required for clustered loci to arise [58]. Both an extension of models predicting genetic

285    architecture and additional empirical studies – on a wider variety organisms and traits – are needed to

286    evaluate the generality of particular architectures and to further understand the conditions under which they

287    evolve.

288        We find additional evidence that a large-effect locus on ssa25, *vgll3,* largely underlies age at maturity

289    in Atlantic salmon corroborating findings from a number of association studies on Atlantic salmon

290    maturation [10,19,31,32,59]. The second strongest associated locus in this study is located in close proximity

291    to *six6* on ssa09. This locus was previously found to be associated with early maturation in male farmed

292    Atlantic salmon [19], with sea age at maturity in wild Atlantic salmon prior to population structure correction

293    [10] and two species of Pacific salmon (Sockeye salmon and Steelhead trout) [36]. Additionally, we found

294    another three loci associated with sea age at maturity: *pecam1, asap2aa* and *taar13c.* The handful of loci

295    found here suggests that wild Atlantic salmon have a mixed genetic architecture where multiple loci, with a

296    variety of effect sizes, control maturation timing – similar to what has been found in male farmed Atlantic

297    salmon [19]. Knowledge of this mixed genetic architecture is highly relevant for how we predict the

298    evolution of maturation timing in wild Atlantic salmon populations. A large body of work has shown the

299    relevance of genetic architecture in determining evolutionary responses [60–68]. Recent works highlight the

300    relevance of the genetic architecture underlying fitness traits when predicting a population's response to

301 environmental changes [69] and selective pressures such a fishing [70]. Future work elucidating how such

302 mixed genetic architectures affect predicted evolution of traits, compared to that of omnigenic or polygenic

303 architectures, will be valuable.

304     We find differences in locations of top-associated SNPs identified here and in Sinclair-Waters et al.

305 [19]. This is not surprising given that we are examining sequence data that captures more SNP variation

306 compared to SNP-array data used in Sinclair-Waters et al. [19]. Furthermore, we failed to find associations

307 between sea age at maturity and many of the candidate regions identified in Sinclair-Waters et al. [19]. For

308 example, several candidate regions on ssa03 and ssa04 displayed particularly strong association signals in

309 aquaculture salmon, however, no signals at these regions were found here. Additionally, only one association

310 peak at ssa06:27541960-28218141 was found here, whereas two independent associations within this region

311 were found in aquaculture salmon [19]. Such differences may reflect changes in the genetic architecture of

312 the trait evolving since the domestication of Atlantic salmon. Although, we would not expect large changes

313 to occur given the domestication is relatively recent, just 10 to 15 generations ago [71]. Furthermore, this

314 study is likely under-powered to detect all previously identified loci, particularly those with smaller effect

315 sizes or rare alleles, due to smaller sample size. Additionally, there could be differences in genetic

316 architecture among environments [72] and/or genotype by environment interactions giving rise to distinct

317 genetic architectures in wild populations versus aquaculture strains.

318     We do not find strong evidence of multi-SNP associations at candidate loci examined in this study,

319 however, we cannot yet disregard the utility of multi-SNP association methods for further resolving the

320 genetic architecture of Atlantic salmon maturation. First, we do not examine the entire genome due to

321 computational restrictions, rather, we focussed on 116 previously identified candidate regions. Second, the

322 Atlantic salmon genome is highly complex [39] and therefore errors in the assembly that may be disruptive

323 for haplotype-based analysis could exist. As new and improved versions of the Atlantic salmon genome are

324 published, our ability to test for haplotypic associations will improve. Furthermore, in a few cases

325 (ssa09:10915066-11415066, ssa09:24636574-25136574, ssa25:28389273-28889273) the *PiMASS* analyses

326 post-regression of the top SNP selected no SNPs for a model explaining sea age at maturity variation,

327 however, the median proportion of variance explained across all iterations was greater than zero. This may

328     suggest that a weak signal was present, but was being missed due to insufficient power. Although this is

329     largely speculative, it suggests that ruling out the possibility of multi-SNP associations at these particular

330     candidate regions may be premature. Higher-powered studies (i.e. more individuals per population) may help

331     to resolve this in the future.

332        In conclusion, our analytical framework, combining both single and multi-SNP association methods,

333     reveals that single SNP variation is sufficient for explaining the association of previously identified

334     candidate loci for Atlantic salmon maturation timing. Previous empirical and theoretical work have described

335     trait-associated loci that have complex alleles with multiple variants, our findings therefore demonstrate the

336     diversity of genetic architectures for fitness-related traits. Additional data, and a greater diversity of species

337     and traits, will serve to better understand why this diversity of genetic architectures exists and how these

338     particular genetic architectures evolve. The analytical framework used here will be a valuable resource for

339     accomplishing this as individual-level resequencing data for wild species with phenotyped individuals

340     becomes increasingly available.

341

353    Research Data (NIRD, project NS9055K). Phenotype data was provided by the Norwegian Institute for

354    Nature Research (NINA).

**Data availability**

356    Genome re-sequencing data for individuals used in this study are available in the European Nucleotide

357    Archive (ENA) or NCBI with the project accession code PRJEB38061 [38].

**Contributions**

359    CRP, NJB, MSW conceived the study. TN developed the variant calling workflow and constructed the fixed

360    assembly of *ssa25*. JW developed the variant filtering criteria. MSW performed all downstream analyses

361    with input from NJB. MPK played key role in generating whole genome sequencing data. SL led the whole

362    genome sequencing work as part of the AquaGenome project. HS, GHB, BFL, CRP coordinated Atlantic

363    salmon sampling and provided phenotypic information. MSW, CRP, NJB drafted the manuscript. All authors

364    commented on and approved the final manuscript.

**Competing interests**

366    There are no competing interests.

367

368    References

369    1.    Marouli E, Graff M, Medina-Gomez C, Lo KS, Wood AR, Kjaer TR, et al. Rare and low-frequency
370          coding variants alter human adult height. Nature. 2017;542(7640):186–90.

371    2.    Timpson NJ, Greenwood CMT, Soranzo N, Lawson DJ, Richards JB. Genetic architecture: the shape
372          of the genetic contribution to human traits and disease. Nat Rev Genet. 2017;

373    3.    Boyle EA, Li YI, Pritchard JK. An Expanded View of Complex Traits: From Polygenic to
374          Omnigenic. Cell. 2017;169(7):1177–86.

375    4.    Moser G, Lee SH, Hayes BJ, Goddard ME, Wray NR, Visscher PM. Simultaneous Discovery,
376          Estimation and Prediction Analysis of Complex Traits Using a Bayesian Mixture Model. PLoS
377          Genet. 2015;11(4):1–22.

378    5.    Loh PR, Bhatia G, Gusev A, Finucane HK, Bulik-Sullivan BK, Pollack SJ, et al. Contrasting genetic
379          architectures of schizophrenia and other complex diseases using fast variance-components analysis.
380          Nat Genet. 2015;47(12):1385–92.

381    6.    Fisher R. The correlations between relatives on the supposition of mendelian inheritance. Philos
382          Trans R Soc Edinburgh. 1918;52:399–433.

383    7.    Pritchard JK, Di Rienzo A. Adaptation - not by sweeps alone. Nat Rev Genet. 2010;11(10):665–7.

384    8.    Liu X, Li YI, Pritchard JK. Trans Effects on Gene Expression Can Drive Omnigenic Inheritance.
385          Cell. 2019;177(4):1022-1034.e6.

386    9.    Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, et al. 10 Years of GWAS
387          Discovery: Biology, Function, and Translation. Am J Hum Genet [Internet]. 2017;101(1):5–22.
388          Available from: http://dx.doi.org/10.1016/j.ajhg.2017.06.005

389    10.   Barson NJ, Aykanat T, Hindar K, Baranski M, Bolstad GH, Fiske P, et al. Sex-dependent dominance
390          at a single locus maintains variation in age at maturity in salmon. Nature. 2015 Dec
391          17;528(7582):405–8.

392    11.   Linnen CR, Poh Y-P, Peterson BK, Barrett RDH, Larson JG, Jensen JD, et al. Adaptive evolution of
393          multiple traits through multiple mutations at a single gene. Science (80- ). 2013;339(6125):1312–6.

394    12.   Lamichhaney S, Fan G, Widemo F, Gunnarsson U, Thalmann DS, Hoeppner MP, et al. Structural
395          genomic changes underlie alternative reproductive strategies in the ruff (*Philomachus pugnax*). Nat
396          Genet. 2015;48(1):84–8.

397    13.   Küpper C, Stocks M, Risse JE, Dos Remedios N, Farrell LL, McRae SB, et al. A supergene
398          determines highly divergent male reproductive morphs in the ruff. Nat Genet. 2015;48(1):79–83.

399    14.   Kirubakaran TG, Grove H, Kent MP, Sandve SR, Baranski M, Nome T, et al. Two adjacent
400          inversions maintain genomic differentiation between migratory and stationary ecotypes of Atlantic
401          cod. Mol Ecol. 2016;25:2130–43.

402    15.   Sinclair-Waters M, Bradbury IR, Morris CJ, Lien S, Kent MP, Bentzen P. Ancient chromosomal
403          rearrangement associated with local adaptation of a post-glacially colonized population of Atlantic
404          Cod in the northwest Atlantic. Mol Ecol [Internet]. 2017;(October):1–13. Available from:
405          http://doi.wiley.com/10.1111/mec.14442

406    16.   Pearse DE, Barson NJ, Nome T, Gao G, Campbell MA, Abadía-Cardoso A, et al. Sex-dependent
407          dominance maintains migration supergene in rainbow trout. bioRxiv [Internet]. 2018;504621.
408          Available from: https://www.biorxiv.org/content/early/2018/12/22/504621.article-metrics

409    17.   Joron M, Frezal L, Jones RT, Chamberlain NL, Lee SF, Haag CR, et al. polymorphic supergene
410          controlling butterfly mimicry. Nature. 2011;

411  18.  Sinnott-Armstrong N, Naqvi S, Rivas MA, Pritchard JK. GWAS of three molecular traits highlights
412       core genes and pathways alongside a highly polygenic background. bioRxiv.
413       2020;2020.04.20.051631.

414  19.  Sinclair-Waters M, Ødegård J, Korsvoll SA, Moen T, Lien S, Primmer CR, et al. Beyond large-effect
415       loci: large-scale GWAS reveals a mixed large-effect and polygenic architecture for age at maturity of
416       Atlantic salmon. Genet Sel Evol [Internet]. 2020;52(1):9. Available from:
417       https://doi.org/10.1186/s12711-020-0529-8

418  20.  Bickel RD, Kopp A, Nuzhdin S V. Composite effects of polymorphisms near multiple regulatory
419       elements create a major-effect QTL. PLoS Genet. 2011;7(1):1–8.

420  21.  Kerdaffrec E, Filiault DL, Korte A, Sasaki E, Nizhynska V, Seren Ü, et al. Multiple alleles at a single
421       locus control seed dormancy in Swedish Arabidopsis. Elife [Internet]. 2016 Dec 14;5(3):1–24.
422       Available from: http://elifesciences.org/lookup/doi/10.7554/eLife.22502

423  22.  Yeaman S, Whitlock MC. The genetic architecture of adaptation under migration-selection balance.
424       Evolution (N Y). 2011;65(7):1897–911.

425  23.  Yeaman S. Genomic rearrangements and the evolution of clusters of locally adaptive loci. Proc Natl
426       Acad Sci U S A [Internet]. 2013;110:E1743-51. Available from:
427       http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3651494&tool=pmcentrez&rendertype=a
428       bstract

429  24.  Remington DL. Alleles versus mutations: Understanding the evolution of genetic architecture
430       requires a molecular perspective on allelic origins. Evolution (N Y). 2015;69(12):3025–38.

431  25.  Stearns SC. Life history evolution: Successes, limitations, and prospects. Naturwissenschaften.
432       2000;87(11):476–86.

433  26.  Mobley KB, Aykanat T, Czorlich Y, House A, Kurko J, Miettinen A, et al. Maturation in Atlantic
434       salmon (Salmo salar, Salmonidae): a review of ecological, genetic, and molecular processes. FEBS
435       Lett. 2020;(November):1–58.

436  27.  Fleming IA, Einum S. Reproductive ecology: a tale of two sexes. In: Atlantic Salmon Ecology. 2011.
437       p. 35–65.

438  28.  Mobley KB, Granroth-Wilding H, Ellmén M, Orell P, Erkinaro J, Primmer CR. Time spent in distinct
439       life history stages has sex-specific effects on reproductive fitness in wild Atlantic salmon. Mol Ecol.
440       2020;29(6):1173–84.

441  29.  Gjerde B. Response to individual selection for age at sexual maturity in Atlantic salmon.
442       Aquaculture. 1984;38(3):229–40.

443  30.  Reed TE, Prodöhl PA, Bradley C, Gilbey J, McGinnity P, Primmer CR, et al. Heritability estimation
444       via molecular pedigree reconstruction in a wild fish population reveals substantial evolutionary
445       potential for sea-age at maturity, but not size within age-classes. Can J Fish Aquat Sci [Internet].
446       2018;cjfas-2018-0123. Available from: http://www.nrcresearchpress.com/doi/10.1139/cjfas-2018-
447       0123

448  31.  Ayllon F, Kjærner-Semb E, Furmanek T, Wennevik V, Solberg MF, Dahle G, et al. The vgll3 Locus
449       Controls Age at Maturity in Wild and Domesticated Atlantic Salmon (Salmo salar L.) Males. PLoS
450       Genet. 2015;11(11):1–15.

451  32.  Sinclair-Waters M, Piavchenko N, Ruokolainen A, Aykanat T, Erkinaro J, Primmer CR. Refining the
452       genomic location of SNP variation affecting Atlantic salmon maturation timing at a key large-effect
453       locus. bioRxiv [Internet]. 2021; Available from:
454       https://www.biorxiv.org/content/early/2021/04/26/2021.04.26.441431

455  33.  Halperin DS, Pan C, Lusis AJ, Tontonoz P. Vestigial-like 3 is an inhibitor of adipocyte

456    differentiation. J Lipid Res. 2013;54(2):473–81.

457   34.   Day FR, Thompson DJ, Helgason H, Chasman DI, Finucane H, Sulem P, et al. Genomic analyses
458          identify hundreds of variants associated with age at menarche and support a role for puberty timing in
459          cancer risk. Nat Genet. 2017;49(6):834–41.

460   35.   Perry JRB, Day F, Elks CE, Sulem P, Thompson DJ, Ferreira T, et al. Parent-of-origin-specific allelic
461          associations among 106 genomic loci for age at menarche. Nature. 2014 Jul 23;514:92.

462   36.   Waters CD, Clemento A, Aykanat T, Garza JC, Naish KA, Narum S, et al. Heterogeneous genetic
463          basis of age at maturity in salmonid fishes. Molcular Ecol. 2021;

464   37.   Cánovas A, Reverter A, DeAtley KL, Ashley RL, Colgrave ML, Fortes MRS, et al. Multi-tissue
465          omics analyses reveal molecular regulatory networks for puberty in composite beef cattle. PLoS One.
466          2014;9(7):1–17.

467   38.   Bertolotti AC, Layer RM, Gundappa MK, Gallagher MD, Pehlivanoglu E, Nome T, et al. The
468          structural variation landscape in 492 Atlantic salmon genomes. Nat Commun [Internet].
469          2020;11(5176). Available from: https://doi.org/10.1101/2020.05.16.099614

470   39.   Lien S, Koop BF, Sandve SR, Miller JR, Kent MP, Nome T, et al. The Atlantic salmon genome
471          provides insights into rediploidization. Nature. 2016 May 12;533(7602):200–5.

472   40.   Chapman B, Kirchner R, Pantano L, Smet M De, Beltrame L, Khotiainsteva T, et al. bcbio/bcbio-
473          nextgen: v1.2.3. 2020 Apr 7 [cited 2020 Sep 17]; Available from: https://zenodo.org/record/3743344

474   41.   Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv
475          [Internet]. 2013 Mar 16 [cited 2020 Sep 17]; Available from: https://arxiv.org/abs/1303.3997

476   42.   Picard toolkit. Broad Institute, GitHub repository. Broad Institute; 2019.

477   43.   Depristo MA, Banks E, Poplin R, Garimella K V., Maguire JR, Hartl C, et al. A framework for
478          variation discovery and genotyping using next-generation DNA sequencing data. Nat Genet.
479          2011;43(5):491–501.

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

483   45.   Li H. SNPable Regions [Internet]. 2009. Available from:
484          http://lh3lh3.users.sourceforge.net/snpable.shtml

485   46.   Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
486          Bioinformatics. 2009 Jul;25(14):1754–60.

487   47.   Danecek P, Auton A, Abecasis G, Albers C, Banks E, DePristo M. The variant call format and
488          vcftools. Bioinformatics. 2011;27(15):2156–2158.

489   48.   Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira M, Bender D. Plink: A tool set for whole-
490          genome association and population-based linkage analyses. Am J Hum Genet [Internet]. 2007;81.
491          Available from: https://doi.org/10.1086/519795

492   49.   Patterson N, Price AL, Reich D. Population Structure and Eigenanalysis. PLoS Genet. 2006;2(12).

493   50.   Quinlan AR, Hall IM. BEDTools: A flexible suite of utilities for comparing genomic features.
494          Bioinformatics. 2010;26(6):841–2.

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

497   52.   Xu H, Guan Y. Detecting local haplotype sharing and haplotype association. Genetics.
498          2014;197(3):823–38.

499    53.    Guan Y. Detecting structure of haplotypes and local ancestry. Genetics. 2014;196(3):625–42.

500    54.    Chang CC, Chow CC, Tellier LCAM, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK:
501          rising to the challenge of larger and richer datasets. Gigascience. 2015 Feb;4(1):7.

502    55.    Guan Y, Stephens M. Bayesian variable selection regression for genome-wide association studies and
503          other large-scale problems. Ann Appl Stat. 2011;5(3):1780–815.

504    56.    Harold Jeffreys. The Theory of Probability. 2020. 470 p.

505    57.    Allendorf FW, Thorgaard GH. Tetraploidy and the Evolution of the Salmonid Fishes. Springer.
506          Monographs in Evolutionary Biology. Boston; 1984. 55–93 p.

507    58.    Yeaman S, Aeschbacher S, Bürger R. The evolution of genomic islands by increased establishment
508          probability of linked alleles. Mol Ecol [Internet]. 2016 Jun 1;25(11):2542–58. Available from:
509          https://doi.org/10.1111/mec.13611

510    59.    Ayllon F, Solberg MF, Glover KA, Mohammadi F, Kjærner-semb E, Fjelldal PG, et al. The influence
511          of vgll3 genotypes on sea age at maturity is altered in farmed mowi strain Atlantic salmon. BMC
512          Genet. 2019;20(44):1–8.

513    60.    Barton NH, Turelli M. Natural and sexual selection on many loci. Genetics. 1991 Jan;127(1):229–55.

514    61.    Turelli M. Heritable genetic variation via mutation-selection balance: Lerch's zeta meets the
515          abdominal bristle. Theor Popul Biol. 1984 Apr;25(2):138–93.

516    62.    Turelli M, Barton NH. Polygenic Variation Maintained by Balancing Selection: Pleiotropy, Sex-
517          Dependent Allelic Effects and G × E Interactions. Genetics. 2004;166(2):1053–79.

518    63.    Turelli M, Barton NH. Dynamics of polygenic characters under selection. Theor Popul Biol.
519          1990;38(1):1–57.

520    64.    Lande R. The maintenance of genetic variability by mutation in a polygenic character with linked
521          loci. Genet Res. 2009/04/14. 1975;26(3):221–35.

522    65.    Bulmer MG. The genetic variability of polygenic characters under optimizing selection, mutation and
523          drift. Genet Res (Camb). 1972;19(1):17–25.

524    66.    Débarre F, Yeaman S, Guillaume F. Evolution of Quantitative Traits under a Migration-Selection
525          Balance : When Does Skew Matter ?*. Am Nat. 2015;186.

526    67.    Fisher R. The genetical theory of natural selection. Clarendon, Oxford; 1930.

527    68.    Yeaman S. Local Adaptation by Alleles of Small Effect *. Am Nat. 2015;186.

528    69.    Kardos M, Luikart G. The genetic architecture of fitness drives population viability during rapid
529          environmental change. Am Nat. 2021;

530    70.    Oomen RA, Kuparinen A, Hutchings JA. Consequences of Single-Locus and Tightly Linked
531          Genomic Architectures for Evolutionary Responses to Environmental Change. J Hered [Internet].
532          2020;319–32. Available from: https://academic.oup.com/jhered/article/111/4/319/5867197

533    71.    Gjerde B, Gjedrem T. Estimates of phenotypic and genetic parameters for carcass traits in Atlantic
534          salmon and rainbow trout. Aquaculture. 1984;36(1–2):97–110.

535    72.    Yan W, Wang B, Chan E, Mitchell-Olds T. Genetic architecture and adaptation of flowering time
536          among environments. New Phytol [Internet]. 2021 Jan 23;n/a(n/a). Available from:
537          https://doi.org/10.1111/nph.17229

538