

2017-09-23

A microsatellite baseline for genetic stock identification of European Atlantic salmon (*Salmo salar* L.)

Gilbey, J

<http://hdl.handle.net/10026.1/10349>

10.1093/icesjms/fsx184

ICES Journal of Marine Science

Oxford University Press (OUP)

All content in PEARL is protected by copyright law. Author manuscripts are made available in accordance with publisher policies. Please cite only the published version using the details provided on the item record or document. In the absence of an open licence (e.g. Creative Commons), permissions for further reuse of content should be sought from the publisher or author.

1 **A microsatellite baseline for genetic stock identification of European Atlantic salmon**

2 **(*Salmo salar* L.).**

3

4 John Gilbey^{1*}, Jamie Coughlan², Vidar Wennevik³, Paulo Prodöhl⁴, Jamie R. Stevens⁵, Carlos

5 Garcia de Leaniz⁶, Dennis Ensing⁷, Eef Cauwelier¹, Corrine Cherbonnel⁸, Sofia Consuegra^{9¶},

6 Mark W. Coulson^{10†}, Tom F. Cross², Walter Crozier⁷, Eileen Dillane², Jonathan S. Ellis^{5†}, Eva

7 García-Vázquez¹¹, Andrew M. Griffiths⁵, Sigurdur Gudjonsson¹², Kjetil Hindar¹³, Sten

8 Karlsson¹³, David Knox¹, Gonzalo Machado-Schiaffino^{11‡}, Dorte Meldrup¹⁴, Einar Eg

9 Nielsen¹⁴, Kristinn Ólafsson¹⁵, Craig R. Primmer^{16‡}, Sergey Prusov¹⁷, Lee Stradmeyer¹, Juha-

10 Pekka Vähä^{16‡}, Alexei Je. Veselov¹⁸, Lucy M.I. Webster¹⁰⁺, Philip McGinnity^{2Δ} and Eric

11 Verspoor^{19Δ}

12

13 ¹*Marine Scotland Science, Freshwater Fisheries Laboratory, Faskally, Pitlochry PH16 5LB, UK*

14 ²*Aquaculture & Fisheries Development Centre, School of Biological, Earth and Environmental*
15 *Sciences, University College, Cork, Ireland*

16 ³*Institute of Marine Research, PO Box 1870 Nordnes, 5817 Bergen, Norway*

17 ⁴*Institute for Global Food Security, School of Biological Sciences, Queen's University, Belfast*
18 *BT9 7BL, UK*

19 ⁵*Department of Biosciences, Geoffrey Pope Building, University of Exeter, Stocker Road,*
20 *Exeter EX4 4QD, UK*

21 ⁶*Department of Biosciences, Swansea University, Swansea, UK*

22 ⁷*Agri-Food and Biosciences Institute Northern Ireland, Fisheries and Aquatic Ecosystems*
23 *Branch, Newforge Lane, Belfast BT9 5PX, UK*

24 ⁸*GENINDEXE, 6 rue des Sports, 17000 La Rochelle, France*

25 ⁹*IBERS, Aberystwyth University, Aberystwyth, UK*

26 ¹⁰*Rivers and Fisheries Trusts of Scotland (RAFTS), CBC House, 24 Canning Street, Edinburgh,*
27 *EH3 8EG, UK*

28 ¹¹*Departament of Functional Biology, Genetics, Universidad de Oviedo, C/Julian Claveria s/n,*
29 *33006 Oviedo, Spain*

30 ¹²*Marine and Freshwater Research Institute, Skúlagata 4, 101 Reykjavík, Iceland*

31 ¹³*Norwegian Institute for Nature Research (NINA), PO Box 5685 Torgard, 7485 Trondheim,*
32 *Norway*

33 ¹⁴*DTU Aqua, National Institute of Aquatic Resources, Technical University of Denmark,*
34 *Vejlsøve 39, 8600 Silkeborg, Denmark*

35 ¹⁵*Matis ohf., Vinlandsleid 12, 113 Reykjavik, Iceland*

36 ¹⁶*Department of Biology, University of Turku, 20014 Turku, Finland*

37 ¹⁷*Knipovich Polar Research Institute of Marine Fisheries and Oceanography, 6 Knipovich*
38 *Street, Murmansk, 183763, Russia*

39 ¹⁸*Institute of Biology, Karelian Research Institute, Pushkinskaya 11, 10 185610 Petrozavodsk,*
40 *Russia*

41 ¹⁹*Marine Scotland Science, Freshwater Fisheries Laboratory, Faskally, Pitlochry PH16 5LB, UK*
42 *and Rivers and Lochs Institute, University of the Highlands and Islands, Inverness College, 1*
43 *Inverness Campus, Inverness, IV2 5NA, UK*

44

45 **Corresponding Author: tel: +44 1796 472060; fax: 01796 473523; email:*

46 *John.Gilbey@gov.scot*

47 ^Δ*Philip McGinnity and Eric Verspoor equally share senior co-authorship*

48

49 [¶]*Present address: Department of Biosciences, Swansea University, Swansea, UK*

50 [‡]*Present address: Rivers and Lochs Institute, University of the Highlands and Islands,*
51 *Inverness College, 1 Inverness Campus, Inverness, IV2 5NA, UK*

52 [†]*Present address: School of Biological and Marine Sciences, University of Plymouth, Drake*
53 *Circus, Plymouth, PL4 8AA, UK*

54 [‡]*Present address: Department of Biology, University of Konstanz, 78457 Konstanz, Germany*

55 [†]*Present address: Science and Advice for Scottish Agriculture, Roddinglaw Road, Edinburgh,*
56 *EH12 9FJ, UK*

57 [‡]*Present address: Department of Biosciences and, Biotechnology Institute, University of*
58 *Helsinki, 00014, Finland*

59 [‡]*Present address: Water and Environment of Western Uusimaa, POB 51, 08101 Lohja,*
60 *Finland*

61

62 *Running title: Stock identification of European Atlantic salmon*

63 **Abstract**

64 Atlantic salmon (*Salmo salar* L.) populations from different river origins mix in the North Atlantic
65 during the marine life stage. To facilitate marine stock identification, we developed a genetic
66 baseline covering the European component of the species' range excluding the Baltic Sea, from the
67 Russian River Megra in the north-east, the Icelandic Ellidaar in the west, and the Spanish Ulla in the
68 south, spanning 3737 km North to South and 2717 km East to West. The baseline encompasses data
69 for 14 microsatellites for 26,822 individual fish from 13 countries, 282 rivers and 467 sampling sites.
70 A hierarchy of regional genetic assignment units was defined using a combination of distance-based
71 and Bayesian clustering. At the top level three assignment units were identified comprising
72 Northern, Southern, and Icelandic regions. A second assignment level was also defined, comprising
73 eighteen and twenty-nine regional units for accurate individual assignment and mixed stock
74 estimates respectively. The baseline provides the most comprehensive geographical coverage for an
75 Atlantic salmon genetic data-set, and a unique resource for the conservation and management of
76 the species in Europe. It is freely available to researchers to facilitate identification of the natal origin
77 of European salmon.
78

79

80 Key words: Atlantic salmon, genetic stock identification, individual assignment, marine ecology, microsatellites

81

82

83

84 **Introduction**

85

86 Homing of Atlantic salmon to natal rivers, in combination with factors such as founder
87 effects, isolation, selection and genetic drift, and broad scale phylogeographic processes,
88 has resulted significant population structuring at a hierarchy of levels from intra-river to
89 inter-continental (King et al., 2001) and locally adapted populations (Garcia de Leaniz et al.,
90 2007) including variations in marine migratory patterns among populations from different
91 parts of the species range (Webb et al., 2007). However, the full extent of differences in

92 migratory patterns among populations and how this may be changing in response to shifting
93 environmental conditions remains to be resolved (Jonsson et al., 2016).

94 Advancing understanding of population and stock-specific migration, distribution and
95 feeding patterns, and their implications for marine mortality rates, and the impact of
96 climate change, is hampered by a lack of information relating to the marine-phase of the
97 lifecycle (Crozier et al., 2004). This situation makes it difficult to appropriately target actions
98 to mitigate anthropogenic influences on different stock components e.g. the impacts of
99 mixed-stock fisheries and bycatches. Thus a tool that allows the accurate identification of
100 genetically distinct populations and regional entities (MacKenzie et al., 2011) and
101 discrimination of the stock origins of fish in mixed feeding aggregations or during migratory
102 phases would be invaluable in species' and North Atlantic marine ecosystem management.

103 DNA profiling methods for identifying the region or river/tributary of origin of
104 salmonids have advanced over recent decades and are widely applied to Pacific salmon
105 (*Oncorhynchus* spp.) stock management (e.g. Shaklee et al., 1999; Beacham et al., 2004;
106 Beacham et al., 2006; Shedd et al., 2016). Their application to Atlantic salmon stock
107 management has provided valuable insights into stock mixing at several spatial scales,
108 including intercontinental (e.g. North American and European stocks in the West Greenland
109 fishery: Gauthier-Ouellet et al., 2009), regional (e.g. stock composition in Canadian gill-net
110 fisheries: Bradbury et al., 2016) and river level (e.g. population structuring in the River
111 Teno/Tana: Vähä et al., 2016). However, overall, its use has been more limited due to the
112 lack of useful genetic baselines for many parts of the species range.

113 Genetic baselines are available for the western side of the Atlantic (e.g. Bradbury et
114 al., 2015; Sheehan et al., 2010), including a recently developed fine scale range-wide North
115 American microsatellite baseline (Bradbury et al., 2016), that facilitate within-region

116 identification of fish originating from Western Atlantic populations at high geographic
117 resolution. In contrast, only partial baselines have been developed for the eastern side of
118 the Atlantic (e.g. Griffiths et al., 2010; Verspoor et al., 2012; Ensing et al., 2013; Gilbey et al.,
119 2016a; Vähä et al., 2016) and no high-resolution baseline exists for the species' non-Baltic,
120 eastern Atlantic range. Such a baseline would allow a DNA-based approach to the genetic
121 stock identification (GSI) of marine samples from the Eastern Atlantic and, in conjunction
122 with ecological studies, would help to provide a more detailed understanding of variations
123 in the North Atlantic migration and distribution patterns of different European Atlantic
124 salmon stocks. Such insight could improve understanding of the factors conditioning marine
125 mortality, and facilitate the implementation of more effective management programmes
126 (Crozier et al., 2004).

127 Genetic stock identification (GSI) has been carried out using various genetic markers,
128 with early work successfully using allozymes (Koljonen and McKinnell, 1996) and
129 mitochondrial DNA (Moriya et al., 2007) for salmonid species in some contexts, including for
130 Atlantic salmon. However, higher levels of resolution and more widespread application has
131 been subsequently achieved using microsatellite loci and they became the genetic marker
132 most widely used in studies of Atlantic salmon stock differentiation. Even though, more
133 recently, attention has turned to Single Nucleotide Polymorphisms (SNPs), the existing large
134 body of microsatellite data available remains a unique and powerful resource that can be
135 exploited for GSI in Atlantic salmon. However, it also has limitations (reviewed in Moran et
136 al., 2006) related to laboratories using different sets of markers, variations in allele-calling
137 with different size markers or allele-size bins, different screening platforms; differences in
138 chemistry, differences in the fluorophore markers across loci and whether the forward or
139 reverse primer is labelled as well as differences in primer sizes. All of these can result in

140 inconsistent allele-size designations across data sets generated by different laboratories.

141 Nevertheless, evidence from large-scale standardisation projects for salmonid species such
142 as *Oncorhynchus mykiss* (Stephenson et al., 2009) and *O. tshawytscha* (Seeb et al., 2007), as
143 well as Atlantic salmon (e.g. Ellis et al., 2011), indicate these issues can be addressed and
144 comprehensive, large scale integrated genetic baselines constructed (Moran et al., 2006).

145 Described here is a trans-European GSI baseline for Atlantic salmon (excluding Baltic
146 salmon stocks which do not migrate to the North Atlantic) constructed by linking existing
147 national and international microsatellite screening programmes. Baltic salmon populations
148 are excluded from the baseline, as they do not migrate outside the Baltic Sea (Karlsson and
149 Karlstrom, 1994; Torniainen et al., 2013). Data was integrated for a common set of 14
150 microsatellite loci for a geographically representative set of rivers spanning the species'
151 Eastern Atlantic European range from the Russian River Megra in the north-east (66.151 N,
152 41.484 W), to the Icelandic Ellidaar in the west (64.117 N, 21.833 E) and the Spanish Ulla in
153 the south (42.639 N, 8.761 E). Baseline samples encompassed rivers responsible for about
154 ≈85% of wild-salmon production in the study region (based on rod-catch data derived from
155 numerous sources). Existing and new data supplied by partners in a multi-laboratory trans-
156 European consortium were calibrated (Ellis et al., 2011), subjected to stringent quality
157 control and integrated to produce the new baseline. A hierarchical assignment unit
158 approach was used and the baseline resolved into genetically distinctive regional
159 assignment units. Assignment power and accuracy to these units were assessed, using both
160 simulations and test samples, the latter constructed by removing fish from the dataset, to
161 establish the utility of the baseline for regional assignment of marine-phase European origin
162 salmon in the North Atlantic.

163

164 **Methods**

165 **Baseline samples**

166 Samples were collected from 32,888 Atlantic salmon from 551 sites representing 325 rivers
167 in 13 countries across Europe (Denmark, England, Finland [two rivers with outlets in
168 Norway], France, Iceland, Ireland, Northern Ireland, Norway, Russia, Scotland, Spain,
169 Sweden and Wales) (Fig. 1, Table 1, Supplementary data S1 & S2), including the Baltic River
170 Torne to act as a genetic out-group. Sampled sites spanned the species' entire eastern
171 Atlantic range and spanned 3737 km from North to South and 2717 km from East to West.

172 Samples were collected from 1994 to 2010, with the majority collected in 2008–2009.
173 Mainly juvenile fish were sampled, mostly parr and fry, but in some cases tissues from
174 smolts or mature salmon returning to fresh water to spawn were sampled. Numbers
175 sampled at a site ranged from 11 to 300 with a mean of 58, and rivers were characterised by
176 1 to 12 sites, depending largely on river size, with a mean number of sample sites per river
177 of 1.7. Full details of sites are given in the Supplementary material (S1 & S2).

178 **Genotyping**

179 Microsatellite data were obtained from DNA extracted from tissue samples (typically fin
180 clips or scales) screened by a consortium of 11 laboratories located across Europe (Table 1)
181 for 14 of the 15 loci identified by a consortium of researchers and described by Olafsson *et*
182 *al.* (2010). *SsaD486* (King *et al.*, 2005) was excluded from the analysis due to its lack of
183 variation over much of the European range. The panel of 14 loci used here were *SsaF43*
184 (Sanchez *et al.*, 1996), *Ssa14*, *Ssa289* (McConnell *et al.*, 1995), *Ssa171*, *Ssa197*, *Ssa202*
185 (O'Reilly *et al.*, 1996) *SSsp1605*, *SSsp2201*, *SSsp2210*, *SSsp2216*, *SSspG7* (Paterson *et al.*,
186 2004), *SsaD144*, *SsaD157* (King *et al.*, 2005) and *SSsp3016* (unpublished, GenBank number
187 AY37820).

188 PCR conditions, thermocyclers and multiplexes varied across laboratories, as did
189 genotyping platforms, size standards and other chemistry employed. Genotyping details and
190 standardisation of genotype assignments among laboratories appear in Ellis et al. (2011). In
191 summary, two 96-well 'control plates' were prepared (Matis, Iceland) containing template
192 DNA extracted from samples representing the widest coverage of the range of *S. salar* as
193 was practicable and which covered sites from both the Eastern and Western Atlantic. These
194 were subsampled and typed by each laboratory. Genotypes were submitted by each
195 member of the consortium to a single depository (Exeter University) where conversion
196 algorithms and standardised nomenclature were applied. For each locus, lists of allele
197 counts and sizes for each laboratory were aligned and cross-referenced for the sample
198 genotypes in the control plates. Standard allele scores were designated for each locus and
199 size differences between allele lists from each laboratory were determined, which allowed
200 laboratory specific standardisation rules to be defined. It should be noted that using this
201 approach not every possible allele was screened, but the approach did allow the individual
202 microsatellite bin ladders to be defined at each location. It cannot be ruled out therefore
203 that rare alleles or alleles affected by regional indels may be have been missed using such an
204 approach, although the coherence of the reference baseline produced (see below) suggests
205 this is unlikely to have been a major influencing factor.

206 Based on the standardisation rules, all data generated for baseline sites were
207 converted to the standard size ranges and stored in a single bespoke database for further
208 analysis (see Ellis et al., 2011 for full details). Sib-ship analysis among individuals in each
209 sample was investigated using the pedigree-likelihood approach implemented within the
210 program COLONY (Jones and Wang, 2010) and used to exclude all but one fish from each
211 full-sib family in each sample prior to inclusion in the database. Fish with less than 10 loci

212 genotyped were removed from further analysis due to concerns with DNA and genotype
213 quality. Sites with more than half of the loci out of Hardy-Weinberg equilibrium (examined
214 in GENEPOP 4.2.2; Rousset, 2008) (potentially not representative of a single population),
215 those that had less than 70% of fish scored at all loci (potentially poor quality DNA and
216 genotypes), and those consisting of less than 30 individuals after quality control checks
217 listed above (potential failure to provide accurate estimates of allele frequencies), were also
218 removed. We estimated descriptive statistics with GenALEx 6 (Peakall and Smouse, 2006).

219 **Assignment units**

220 Assignment units were defined in an iterative way similar to that employed by Gilbey et al.
221 (2016a). Units were first defined by a combination of distance-based and Bayesian
222 clustering. Individual assignment accuracies using these units were then examined and units
223 where accuracies did not meet a predefined threshold were combined with units that saw
224 reciprocal misassignments, until all units had accuracies at or above the threshold level.

225 The distance-based approach was based on a neighbour-joining tree (Saitou and Nei,
226 1987) constructed using Nei's genetic distance D_A (Nei et al., 1983) calculated in POPTREE2
227 (Takezaki et al., 2010) and visualised in MEGA7 (Kumar et al., 2016). The clustering approach
228 was carried out in STRUCTURE (Pritchard et al., 2000), using a burn-in of 100,000 and a run
229 phase of 300,000 iterations during each application. Three replicates for each cluster
230 number (K) were run with values of K from 1 to 10. $K = 10$ emerged as an upper limit after
231 monitoring of the results of the runs while they were underway. In each case stable
232 estimates of true K at the level under analysis had been identified by this point (see results).
233 Prior site information was incorporated into the analysis using the LOCPRIOR option. The
234 smallest K capturing the major structure in the dataset was defined by the ΔK method of
235 Evanno et al. (2005), which was calculated using STRUCTURE HARVESTER (Earl and

236 vonHoldt, 2012). Replicate membership coefficients were combined with CLUMPP

237 (Jakobsson and Rosenberg, 2007) using the Full Search method.

238 The Bayesian clustering was carried out using a hierarchical approach, starting with
239 the full dataset. Evanno et al. (2005) showed that STRUCTURE tends to capture the major
240 structure in a reference dataset but that more fine scale structure may become evident if a
241 hierarchical analysis is performed. In the current analysis, at each hierarchical level a
242 STRUCTURE analysis was performed and the minimum best K identified. The data were then
243 split up into the cluster units and further STRUCTURE analysis performed on each one
244 independently. This was repeated at each hierarchical split until either single-river
245 structuring was observed or geographical coherence of the clusters was lost.

246 Once both the distance-based and clustering analysis had been performed, the degree
247 to which the assignment units identified by each technique corresponded was examined.
248 Where the same units were identified these were incorporated into the initial assignment
249 unit panel. Where the two approaches had identified different units the smallest unit from
250 either approach was incorporated into the initial assignment unit panel, for example in a
251 situation where one technique had identified a single unit and another had identified sub-
252 units the sub-units were added to the initial assignment panel. In this way, the smallest
253 units identified by one or both technique were incorporated into the initial assignment unit
254 test panel.

255 Once the initial assignment unit panel had been identified, individual assignment
256 accuracy was calculated for each of these units (see below). If the assessed accuracy to a
257 unit was at or above 80% the unit was retained in the panel. If accuracy was below this level
258 the unit was combined with other units to which reciprocal misassignments were occurring.
259 Accuracies were tested again and the process repeated until all units in the panel had

260 individual assignment accuracies at or above the 80% level. Nei's genetic distance D_A (Nei et
261 al., 1983) was again calculated for all pairwise final assignment combinations using the
262 POPULATIONS 1.2.3 software package (Langella, 1999).

263 **Assignment analysis**

264 *Individual assignment*

265 Individual assignment accuracy was calculated using maximum likelihood-based mixture
266 analyses carried out using ONCOR (Kalinowski et al., 2007) with mixture proportions
267 estimated using the EM algorithm and genotype probabilities calculated by the method of
268 Rannala and Mountain (1997). In order to estimate unbiased assignment accuracies using
269 fish not represented in the baseline, assignment tests were based on fish randomly
270 removed from the reference baseline and combined into a mixture file. A randomly selected
271 10% of fish were removed from each of the three top level assignment units identified (see
272 results) resulting in a total of 2682 fish in the mixture file. For each fish the most likely
273 assignment unit of origin and associated assignment probability was calculated. Fish with
274 assignment probabilities below 0.8 were classified as unassigned and excluded from the
275 analysis. Accuracy to the assignment units was then calculated with the remaining fish.
276 Using such a cut-off meant that fish whose origin was difficult to determine (low probability)
277 were removed from the analysis and so potential accuracy could be increased (Gilbey et al.,
278 2016a; Bekkevold et al., 2015). However, the application of cut-off scores also increased the
279 proportion of unassigned fish (Gilbey et al., 2016a) and can thus influence apparent stock
280 proportions if calculated from the individual assignments. As such, this should not be
281 performed for this purpose and so, in order to estimate accurate stock proportions a Mixed
282 Stock Analysis (MSA) approach was utilised (see below).

283 *100% simulations*

284 Simulated fishery mixtures were analysed in ONCOR and comprised sets of 100% simulated
285 samples of fish from each assignment unit. Genotypic frequencies for each locus in each unit
286 were re-sampled following Anderson *et al.* (2008). The 100% simulations were based on
287 1000 simulations of 200 fish per hierarchical assignment unit and the same simulated
288 reference sample sizes as in the actual dataset.

289 *Mixed stock analysis*

290 Mixed stock proportions were calculated for each assignment unit. The same set of 2682
291 randomly selected fish used for the individual assignments was used and mixture
292 proportions estimated in ONCOR using conditional maximum likelihood (Millar, 1987) with
293 confidence intervals calculated based on 1000 bootstraps.

294 *Equal proportions*

295 Mixed stock proportions were calculated for each assignment unit using simulated fishery
296 mixtures with equal proportions of fish at each assignment unit in ONCOR. One hundred fish
297 were simulated for each unit and confidence intervals of the estimates calculated using
298 1000 bootstraps.

299 *Baseline coverage analysis – River removal*

300 A baseline rarely covers all possible source populations completely, and so some fish in real
301 fishery mixtures may be from populations not included in the baseline. Hence, simulation
302 analysis may overestimate the success rates of assignments of fish in an actual fishery due
303 to being based only on samples from sites and rivers contained in the baseline (Waples *et al.*,
304 2008). This issue was addressed using a further test panel and associated test baseline. A
305 random 10% of the rivers in each assignment unit were removed from the baseline and used
306 as test mixtures that were then assigned back to the reconstructed baseline. All assignment
307 units comprising more than one river had at least one river randomly removed (see

308 Supplementary material S1 for details of sites and rivers removed). Fish in these
309 'unrepresented' mixture panels were thus from sites and rivers not included in the
310 reconstructed baseline. In this way, we tested the capability of the baseline to reflect the
311 regional signal of each assignment unit and to assign fish from sites and rivers not included
312 in the baseline but from the assignment unit. This procedure was repeated at both
313 assignment unit levels, again using ONCOR, with confidence intervals calculated based on
314 1000 bootstraps.

315

316 **Results**

317 **Baseline QC**

318 From a total of 551 sites sampled, 84 sites were removed, leaving 467 sites containing
319 26,822 fish representing 282 rivers in the final baseline (Table 1). From those removed, 17
320 sites were not in H-W proportions, 51 had <70% of fish screened at all loci, and 15 had <30
321 individuals representing the site after correction for full-siblings and individual fish for which
322 <10 loci could be reliably genotyped. A further site (a sample of adult rod-caught fish from
323 the Norwegian River Flekkeelva in 2007) was removed due to extreme outlier behaviour in
324 the STRUCTURE analysis (data not shown). Full site details are contained in Fig. 1, Table 1
325 and Supplementary data S1 & S2. Across sites most loci were highly variable, with allele
326 numbers ranging from 10 for *Ssa14* to 46 for *SsaD157* (mean 29.9). Additional descriptive
327 and diversity estimates for each locus and site are presented in Supplementary material S3.

328 **Definition of initial assignment regions**

329 A neighbour-joining tree of Nei's D_A is summarised in Fig. 2 with an expanded version
330 detailed in Supplementary data S4 and full site level D_A matrix in Supplementary data S5. A
331 plot of ΔK , and a map showing the geographic positioning of the clusters at each hierarchical

332 STRUCTURE level are shown in Fig. 3. Assignment units as defined by POPTREE and
333 STRUCTURE are compared in Supplementary data S6.

334 Both distance-based N-J tree and Bayesian STRUCTURE approaches identified three
335 large regional groupings of sites covering the Northern, Southern and Icelandic regions and
336 these will henceforth be referred to as the Level 1 assignment units. There was in general a
337 good agreement between the two population structuring techniques at the lowest level
338 units identified. Indeed, of the 26 and 22 units defined by the NJ Tree and Bayesian
339 clustering methods, respectively, 17 units were identical (Supplementary data S6). Using the
340 lowest level divisions produced from each technique resulted in a total of 29 units identified
341 for the initial Level 2 assignment accuracy testing (column 1 in Table 2, Supplementary data
342 S6). The assignment units at both initial levels are mapped in Fig. 1, with D_A matrixes
343 detailed in Supplementary data S8.

344 **Assignment analysis**

345 *Initial assignment accuracy*

346 Using the 2682 fish removed from the baseline, individual assignments were performed at
347 Level 1 and at the initially defined Level 2 assignment units. At Level 1 the assignment
348 accuracy of all fish to the Northern, Southern and Icelandic unit respectively was 90.8%,
349 92.7% and 99.5% respectively. Using a probability cut-off score ≥ 0.8 this increased to
350 94.2%, 95.5% and 100% with 86.8%, 90.2% and 99.5% of fish in the mixture being assigned.

351 Assignment accuracy of fish with probability scores ≥ 0.8 to the Level 2 units was \geq
352 80% in 19 of the 29 units (Table 2; for full breakdown of assignments at each Level 2
353 iterative level see Supplementary data S7). After combining assignment units based on
354 reciprocal misassignments, 21 assignment units remained with recalculated accuracies \geq

355 80%. A final round of assignment unit combination resulted in 18 assignment units for which
356 assignment accuracies were all $\geq 80\%$ (Table 2, Supplementary data S7).

357 *100% simulations*

358 The 100% simulations for each assignment unit showed robust estimates of stock
359 proportions at both assignment levels (Fig. 4). At Level 1, the mean estimates matched the
360 actual proportions extremely well with a maximum difference of just 0.3% between the
361 actual and estimated values and all upper CI at 100%. The initial Level 2 assignment units
362 again showed relatively accurate estimates with an average difference between the
363 estimated and actual mean proportions of 4.5%. The West and Central Scotland level,
364 however, showed a difference of 17.6% between estimated and actual proportions. At the
365 first round of assignment unit combinations accuracies were seen to improve, as expected,
366 with average and maximum differences between the estimated and actual mean
367 proportions of 4.5% and 9.0%. These levels reduced to 1.9% and 8.0% respectively at the
368 final Level 2 assignment unit combination round.

369 *Mixed stock analysis*

370 The results of the MSA using the 2682 fish removed from the baseline and used as a fishery
371 mixture are shown in Fig. 5A. For all assignment units, within both assignment levels, apart
372 from a single unit in Level 2, South France/Spain, where the upper CI was 0.19 below the
373 actual value, the estimated proportions of fish in the unit mixtures matched actual
374 proportions (i.e. were within the CI bands). The estimates were also very precise with
375 average CI bands of just 2.2 and a maximum of 4.7. Considering the high accuracy of the
376 mixed stock estimates at this initial assignment unit composition, no further assignment unit
377 amalgamations were deemed necessary for mixed stock analysis.

378 *Equal proportions*

379 As with the previous analysis the equal proportion simulation showed excellent agreement
380 between the actual and estimated proportions in the mixture (Fig. 5B). At Level 1 there was
381 an average difference between actual and estimated of just 0.06% and a maximum of 0.09%
382 (Southern unit) and at Level 2 these two differences only rise to a mean difference of 0.4
383 and a maximum of 1.1% (North Ireland unit).

384 *Baseline coverage analysis – River removal*

385 The most demanding test of assignment capabilities of the baseline was the “river removal”
386 test in which entire river systems were removed from the baseline and their fish assigned to
387 region of origin using the remainder of the rivers in the reference baseline. However, even
388 here relatively high levels of assignment accuracy were obtained (Fig. 5C). Average
389 differences between actual and estimated mixture proportions were 1.9% with a maximum
390 of 2.3% (Southern unit) at Level 1 and 1.3% and 2.9% (Central Scotland/North England)
391 respectively at Level 2. At no time were significant proportions assigned to any of the six
392 single-river assignment units which were not represented in the mixture file (lower CI at
393 zero in these units).

394

395 **Discussion**

396 The study, encompassing the largest analysis of Atlantic salmon population structure in the
397 Eastern Atlantic, for the first time, provides a genetic framework to exploit the power of
398 microsatellite variation to assign Atlantic salmon from this part of the species' range to
399 smaller scale regional stock groups. As such, the reported genetic baseline provides a
400 powerful resource that can be used to increase understanding of the biology of European
401 Atlantic salmon stocks in the North Atlantic marine environment. Enhanced understanding
402 of stock-specific marine migration, distribution, feeding patterns, exploitation and mortality

403 rates, will help to provide guidance towards a more efficient management of Atlantic
404 salmon in a changing environment (Crozier et al., 2004).

405 Distance-based and Bayesian cluster based analyses both reveal hierarchical
406 structuring of river populations of European and Icelandic salmon into regional groups. At
407 the highest level, this structure encompasses large-scale geographical discontinuities
408 between northern (Scandinavia-Russia), Icelandic, and southern regions (Britain-Ireland-
409 France-Denmark-Spain). Such differences have been identified in previous analyses of
410 Atlantic salmon population structure. For example, King et al. (2001) showed with
411 microsatellites an unambiguous separation of Iceland, Norway and Scotland-Ireland-Spain
412 (their Fig. 3), and Verspoor et al. (2005) identified an Icelandic group together with a
413 southern British Isles-Northern France group using allozymes, although a more complex
414 pattern was apparent in their analysis among the more central range groups.

415 At the next highest level, two assignment units shared the largest average degree of
416 distinctiveness from other units, the two also being on opposite extremes of the neighbour-
417 joining tree (Fig. 2). The Baltic unit had a mean D_A of 0.236 to other units (Supplementary
418 data S8), a level of differentiation to other European rivers seen in previous studies (Bourret
419 et al., 2013) and consistent with the restricted migration of Baltic stocks (Karlsson and
420 Karlstrom, 1994) and their long history of geographical isolation (Bourret et al., 2013). A
421 second assignment unit, the English Chalk streams, also shared a similarly high mean D_A of
422 0.236. Griffiths et al. (2010) and Ikediashi et al. (2018) also reported these rivers in Southern
423 England to be highly differentiated from others in the southern part of the European range.
424 However, it is unexpected in the context of the entire European and Icelandic range, that
425 the degree of differentiation matches that of the Baltic.

426 Within Iceland the salmon populations segregate into Northern and Western
427 Icelandic units as was also reported by Olafsson et al. (2014) which is thought may reflect
428 the patterns of recolonisation after the Last Glacial Maximum.

429 Initially the Northern Level 2 unit subdivided into eleven geographically coherent
430 genetic clusters that matched well with previously reported structure in this region. Bourret
431 et al. (2013), using SNP markers, found separation of northern Norway and Russian rivers
432 from the Norwegian and Swedish Atlantic coast rivers, and Kjærner-Semb et al. (2016) found
433 separation of northern and southern Norwegian groupings. Within the northern Norway-
434 Russian complex, Ozerov et al. (2017) also found the same North Kola, Northern Norway
435 and Russia-White sea units as reported here. However, their use of 33 microsatellites and a
436 more comprehensive geographical coverage allowed them to define structure at further
437 hierarchical levels within these groups unresolved in the present study using only 14
438 microsatellites and more limited population coverage.

439 The population structuring of rivers from across the part of the range covered by the
440 Level 1 Southern unit into an initial sixteen Level 2 units accords well with that reported by
441 Griffiths et al. (2010) based on 12 microsatellites, 11 of which form part of in the panel used
442 in the present study. Their study encompassed fish from 57 rivers across the Southern
443 region but excluded rivers from the East coast of Scotland and Northern Ireland and showed
444 similar geographic patterns of genetic structure (their Fig. 2). Similar assignment units in
445 France and Northern Spain appeared in both analyses and also broadly reflected allozyme-
446 based regional differentiation (Verspoor et al., 2005). However, some differences were seen
447 with some of the units between the two methods used to resolve assignment units. Griffiths
448 et al. (2010) identified groupings stretching across both Scotland and Ireland (see their Fig.
449 2) and similar groups were identified here using the STRUCTURE based approach (Fig. 3). In

450 contrast, using the distance-based approach the various Scottish and Irish units were clearly
451 separated (Fig. 2) to which generally good assignments of fish could be made. Nevertheless,
452 some reciprocal misassignment was still evident (Supplementary data S7) suggesting a
453 degree of homology between the units. Further, finer-scale investigation is perhaps required
454 to disentangle completely the complex patterns of population grouping within these
455 regions.

456 Accurate assignments to the initial Level 2 units was not possible at the individual
457 level but was achieved for mixed stock fishery estimates. Acceptable levels of individual
458 assignments could be made to some defined units using the initial split but some areas
459 proved problematic at this scale particularly for Britain and Ireland. This difference reflects
460 the differing power of the two IA and MSA techniques (Manel et al., 2005) and suggests
461 that, when using the baseline for a particular purpose, the required levels of both accuracy
462 and resolution should be defined *a priori*. In turn, this will depend on the specific questions
463 being examined and the tools being utilised.

464 Overall, the two levels of genetic structure are geographically consistent and in basic
465 agreement with major regional phylogeographic groups previously reported using a variety
466 of markers, suggesting the higher level regional structuring is geographically and temporally
467 robust. In contrast, differentiation between regional units identified at the finer geographic
468 scales may in part be conditioned by human activities, such as the transport and escape of
469 fish from aquaculture facilities, stocking, habitat alteration, fisheries-induced evolution, and
470 indirect genetic changes from disease and ecological disturbances. Such genetic structuring,
471 if defined by such contemporary influences, may not have temporal stability and such lower
472 level units thus will need to be monitored to determine if they are stable. Encouragingly, in
473 a previous assessment of temporal stability on assignment of Atlantic salmon in the species'

474 southern European range (Griffiths et al., 2010), test samples collected 20 years before the
475 baseline samples still showed predominant allocation back to region of origin. This finding
476 suggests, at least at the larger scale, regional level units are likely to be temporarily stable.
477 However, this should not be assumed to always be the case and a program of resampling
478 should be incorporated if the baseline is exploited in the future.

479 For the Level 1 and the final Level 2 regional units, all tests of power suggest high
480 accuracies can be achieved with both individual assignments and mixed stock analysis.
481 Accuracies are improved by use of a probability cut-off of 0.8 for individual assignments,
482 which may be useful in some contexts. However, this will reduce the proportion of fish
483 assigned. Thus in application, the best cut-off will depend on the question address and will
484 need to be decided by each individual user. This will also apply to the assignment units used;
485 if reduced accuracies to some of the combined units are acceptable these may also be used
486 in specific circumstances.

487 The assignment tests carried out indicate that the described baseline can be
488 exploited to help investigate patterns of ocean utilisation and associated differences in
489 marine mortality operating at the regional stock level. However, important quantitative
490 variation linked to how individual population components use the ocean, which may affect
491 mortality rates, also exists at the level of individual rivers within regions and among river
492 tributaries (Barson et al., 2015). Evaluation of river-specific problems, likely to exist in some
493 contexts, will require assignments at the individual river level, for which the current baseline
494 appears to have limited usefulness. Nevertheless, even if river-level identification is
495 problematic, identification of region of origin may allow finer scale analysis using higher
496 resolution region-specific baselines.

497 Resolution of intra-regional population contributions in mixed oceanic samples,
498 including within-river contribution assessments, would be facilitated by further increases in
499 the coverage and resolution of the baseline. Higher resolution is already being achieved in
500 selected areas covered by the baseline reported here (Gilbey et al., 2016a; Ozerov et al.,
501 2017; Vähä et al., 2016). Ideally, future work will likely increase baseline coverage to include
502 most of the estimated 2000 rivers in the North-East Atlantic Commission area. However, this
503 will involve diminishing returns given that the rivers currently in the baseline represent an
504 estimated ≈85% of the non-Baltic European adult salmon production. Nevertheless, genetic
505 characterisation of as many populations as possible will be important for biodiversity
506 inventory and assessment. Considerable value could also be added by combining the
507 European baseline reported here with North American information to provide a trans-ocean
508 baseline and thus enable oceanic scale investigations. This has already started using a
509 reduced set of microsatellite markers and shows promise in the ability to assign fish from
510 the entire species' range (Gilbey et al., 2016b).

511

512 **Supplementary data**

513 Supplementary material

514 is available at the ICESJMS online version of the manuscript.

515

516 **Acknowledgments**

517 This work forms part of the SALSEA-Merge research project (Project No. 212529) and was

518 funded by the European Union under theme six of the 7th Framework programme. It was

519 also co-sponsored by the Atlantic Salmon Trust and the Total Foundation, who we thank for

520 financial support. PMcG and JC were partly supported by the Beaufort Marine Research
521 Award in Fish Population Genetics funded by the Irish Government under the Sea Change
522 Programme. The work was also supported under financial support of the program of
523 fundamental research of Presidium of RAS "Searching fundamental scientific investigations
524 in the interests of development of the Arctic zone of Russian Federation". The authors also
525 thank the numerous people responsible across Europe whom helped collect samples and
526 assisted with the laboratory analyses. The manuscript benefited greatly from editorial and
527 reviewer comments on an earlier draft and the authors wish to express thanks for their time
528 and comprehensive inputs.

529

530 **References**

531

- 532 Anderson, E. C., Waples, R. S., and Kalinowski, S. T. 2008. An improved method for predicting the
533 accuracy of genetic stock identification. *Canadian Journal of Fisheries and Aquatic Sciences*,
534 65: 1475-1486.
- 535 Barson, N. J., Aykanat, T., Hindar, K., Baranski, M., Bolstad, G. H., Fiske, P., Jacq, C., et al. 2015. Sex-
536 dependent dominance at a single locus maintains variation in age at maturity in salmon.
537 *Nature*, 528: 405-408.
- 538 Beacham, T., Lapointe, M., Candy, J., Miller, K., and Withler, R. 2004. DNA in action: rapid application
539 of DNA variation to sockeye salmon fisheries management. *Conservation Genetics*, 5: 411-
540 416.
- 541 Beacham, T. D., Candy, J. R., Jonsen, K. L., Supernault, J., Wetklo, M., Deng, L., Miller, K. M., et al.
542 2006. Estimation of Stock Composition and Individual Identification of Chinook Salmon
543 across the Pacific Rim by Use of Microsatellite Variation. *Transactions of the American*
544 *Fisheries Society*, 135: 861-888.
- 545 Bekkevold, D., Helyar, S. J., Limborg, M. T., Nielsen, E. E., Hemmer-Hansen, J., Clausen, L. A. W., and
546 Carvalho, G. R. 2015. Gene-associated markers can assign origin in a weakly structured fish,
547 Atlantic herring. *ICES Journal of Marine Science*, 72: 1790-1801.
- 548 Bourret, V., Kent, M. P., Primmer, C. R., Vasemägi, A., Karlsson, S., Hindar, K., McGinnity, P., et al.
549 2013. SNP-array reveals genome-wide patterns of geographical and potential adaptive
550 divergence across the natural range of Atlantic salmon (*Salmo salar*). *Molecular Ecology*, 22:
551 19.
- 552 Bradbury, I. R., Hamilton, L. C., Chaput, G., Robertson, M. J., Goraguer, H., Walsh, A., Morris, V., et al.
553 2016. Genetic mixed stock analysis of an interceptory Atlantic salmon fishery in the
554 Northwest Atlantic. *Fisheries Research*, 174: 234-244.
- 555 Bradbury, I. R., Hamilton, L. C., Rafferty, S., Meerburg, D., Poole, R., Dempson, J. B., Robertson, M. J.,
556 et al. 2015. Genetic evidence of local exploitation of Atlantic salmon in a coastal subsistence

- 557 fishery in the Northwest Atlantic. *Canadian Journal of Fisheries and Aquatic Sciences*, 72: 83-
558 95.
- 559 Crozier, W. W., Schön, P.-J., Chaput, G., Potter, E. C. E., Maoiléidigh, N. Ó., and MacLean, J. C. 2004.
560 Managing Atlantic salmon (*Salmo salar* L.) in the mixed stock environment: challenges and
561 considerations. *ICES Journal of Marine Science*, 61: 1344-1358.
- 562 Earl, D., and vonHoldt, B. 2012. STRUCTURE HARVESTER: a website and program for visualizing
563 STRUCTURE output and implementing the Evanno method. *Conservation Genetics*
564 *Resources*, 4: 359-361.
- 565 Ellis, J. S., Gilbey, J., Armstrong, A., Balstad, T., Cauwelier, E., Cherbonnel, C., Consuegra, S., et al.
566 2011. Microsatellite standardization and evaluation of genotyping error in a large multi-
567 partner research programme for conservation of Atlantic salmon (*Salmo salar* L.). *Genetica*,
568 139: 353-367.
- 569 Ensing, D., Crozier, W. W., Boylan, P., O'Maoiléidigh, N., and McGinnity, P. 2013. An analysis of
570 genetic stock identification on a small geographical scale using microsatellite markers, and
571 its application in the management of a mixed-stock fishery for Atlantic salmon *Salmo salar* in
572 Ireland. *Journal of Fish Biology*, 82: 2080-2094.
- 573 Evanno, G., Regnaut, S., and Goudet, J. 2005. Detecting the number of clusters of individuals using
574 the software structure: a simulation study. *Molecular Ecology*, 14: 2611-2620.
- 575 Garcia de Leaniz, C., Fleming, I. A., Einum, S., Verspoor, E., Jordan, W. C., Consuegra, S., Aubin-Horth,
576 N., et al. 2007. A critical review of adaptive genetic variation in Atlantic salmon: implications
577 for conservation. *Biological Reviews*, 82: 173-211.
- 578 Gauthier-Ouellet, M., Dionne, M., Ianie, Caron, F., ois, King, T. L., and Bernatchez, L. 2009.
579 Spatiotemporal dynamics of the Atlantic salmon (*Salmo salar*) Greenland fishery inferred
580 from mixed-stock analysis. *Canadian Journal of Fisheries and Aquatic Sciences*, 66: 2040-
581 2051.
- 582 Gilbey, J., Cauwelier, E., Coulson, M. W., Stradmeyer, L., Sampayo, J. N., Armstrong, A., Verspoor, E.,
583 et al. 2016a. Accuracy of Assignment of Atlantic Salmon (*Salmo salar* L.) to Rivers and
584 Regions in Scotland and Northeast England Based on Single Nucleotide Polymorphism (SNP)
585 Markers. *PLoS ONE*, 11: e0164327.
- 586 Gilbey, J., Wennevik, V., Bradbury, I. R., Fiske, P., P., H. L., Jacobsen, J. A., and Potter, T. 2016b.
587 Genetic stock identification of Atlantic salmon caught in the Faroes fishery. *Fisheries*
588 *Research*, 187: 110-119.
- 589 Griffiths, A. M., Machado-Schiaffino, G., Dillane, E., Coughlan, J., Horreo, J. L., Bowkett, A. E.,
590 Minting, P., et al. 2010. Genetic stock identification of Atlantic salmon (*Salmo salar*)
591 populations in the southern part of the European range. *BMC Genetics*, 11:31.
- 592 Ikediashi, C. I., Paris, J. R., King, R. A., Ibbotson, A., and Stevens, J. R. 2018. Atlantic salmon (*Salmo*
593 *salar* L.) in the chalk streams of England are genetically unique. *Journal of Fish Biology*, 92: In
594 Press.
- 595 Jakobsson, M., and Rosenberg, N. A. 2007. CLUMPP: a cluster matching and permutation program
596 for dealing with label switching and multimodality in analysis of population structure.
597 *Bioinformatics*, 23: 1801-1806.
- 598 Jones, O. R., and Wang, J. 2010. COLONY: a program for parentage and sibship inference from
599 multilocus genotype data. *Molecular Ecology Resources*, 10: 551-555.
- 600 Jonsson, B., Jonsson, N., and Albrechtsen, J. 2016. Environmental change influences the life history of
601 salmon *Salmo salar* in the North Atlantic Ocean. *Journal of Fish Biology*, 88: 618-637.
- 602 Kalinowski, S. T., Manlove, K. R., and Taper, M. L. 2007. ONCOR: a computer program for genetic stock
603 identification. Department of Ecology, Montana State University. Available from
604 <http://www.montana.edu/kalinowski/Software/ONCOR.htm>.
- 605 Karlsson, L., and Karlstrom, O. 1994. The Baltic salmon (*Salmo salar* L.): its history, present situation
606 and future. *Dana*, 10: 24.

- 607 King, T., Eackles, M., and Letcher, B. 2005. Microsatellite DNA markers for the study of Atlantic
608 salmon (*Salmo salar*) kinship, population structure, and mixed-fishery analyses. *Molecular*
609 *Ecology Notes*, 5: 130-132.
- 610 King, T. L., Kalinowski, S. T., Schill, W. B., Spidle, A. P., and Lubinski, B. A. 2001. Population structure
611 of Atlantic salmon (*Salmo salar* L.): a range-wide perspective from microsatellite DNA
612 variation. *Molecular Ecology*, 10: 807-821.
- 613 Kjaerner-Semb, E., Ayllon, F., Furmanek, T., Wennevik, V., Dahle, G., Niemela, E., Ozerov, M., et al.
614 2016. Atlantic salmon populations reveal adaptive divergence of immune related genes - a
615 duplicated genome under selection. *BMC Genomics*, 17: 610.
- 616 Koljonen, M. L., and McKinnell, S. 1996. Assessing seasonal changes in stock composition of Atlantic
617 salmon catches in the Baltic Sea with genetic stock identification. *Journal of Fish Biology*, 49:
618 998-1018.
- 619 Kumar, S., Stecher, G., and Tamura, K. 2016. MEGA7: Molecular Evolutionary Genetics Analysis
620 version 7.0 for bigger datasets. *Molecular Biology and Evolution*, 33: 1870-1874.
- 621 Langella, O. 1999. Populations 1.2.32: a population genetic software. CNRS UPR9034. Available at
622 <http://bioinformatics.org/~tryphon/populations/>.
- 623 MacKenzie, K. M., Palmer, M. R., Moore, A., Ibbotson, A. T., Beaumont, W. R. C., Poulter, D. J. S., and
624 Trueman, C. N. 2011. Locations of marine animals revealed by carbon isotopes. *Scientific*
625 *Reports*, 1: DOI: 10.1038/srep00021.
- 626 Manel, S., Gaggiotti, O. E., and Waples, R. S. 2005. Assignment methods: matching biological
627 questions with appropriate techniques. *Trends in Ecology and Evolution*, 20: 136-142.
- 628 McConnell, S. K., Hamilton, L., Morris, D., Cook, D., Paquet, D., and Bentzen, P. 1995. Isolation of
629 salmonid microsatellite loci and their application to the population genetics of Canadian
630 stocks of Atlantic salmon. *Aquaculture*, 137: 19-30.
- 631 Millar, R. B. 1987. Maximum likelihood estimation of mixed stock fishery composition. *Canadian*
632 *Journal of Fisheries and Aquatic Sciences*, 44: 583-590.
- 633 Moran, P., Teel, D. J., LaHood, E. S., Drake, J., and Kalinowski, S. 2006. Standardising multi-laboratory
634 microsatellite data in Pacific salmon: an historical view of the future. *Ecology of Freshwater*
635 *Fish*, 15: 597-605.
- 636 Moriya, S., Sato, S., Azumaya, T., Suzuki, O., Urawa, S., Urano, A., and Abe, S. 2007. Genetic Stock
637 Identification of Chum Salmon in the Bering Sea and North Pacific Ocean Using
638 Mitochondrial DNA Microarray. *Marine Biotechnology*, 9: 179-191.
- 639 Nei, M., Tajima, F., and Tatenno, Y. 1983. Accuracy of estimated phylogenetic trees from molecular
640 data. *Journal of Molecular Evolution*, 19: 153-170.
- 641 O'Reilly, P. T., Hamilton, L. C., McConnell, S. K., and Wright, J. M. 1996. Rapid analysis of genetic
642 variation in Atlantic salmon (*Salmo salar*) by PCR multiplexing of dinucleotide and
643 tetranucleotide microsatellites. *Canadian Journal of Fisheries and Aquatic Sciences*, 53:
644 2292-2298.
- 645 Olafsson, K., Hjorleifsdottir, S., Pampoulie, C., Hreggvidsson, G. O., and Gudjonsson, S. 2010. Novel
646 set of multiplex assay (SalPrint15) for efficient analysis of 15 microsatellite loci of
647 contemporary samples of the Atlantic salmon (*Salmo salar*). *Molecular Ecology Resources*,
648 10: 533-537.
- 649 Olafsson, K., Pampoulie, C., Hjorleifsdottir, S., Gudjonsson, S., and Hreggvidsson, G. O. 2014. Present-
650 Day Genetic Structure of Atlantic Salmon (*Salmo salar*) in Icelandic Rivers and Ice-Cap
651 Retreat Models. *PLoS ONE*, 9: e86809.
- 652 Ozerov, M., Vähä, J. P., Wennevik, V., Niemelä, E., Svenning, M., Prusov, S., Diaz Fernandez, R., et al.
653 2017. Comprehensive microsatellite baseline for genetic stock identification of Atlantic
654 salmon (*Salmo salar* L.) in northernmost Europe. *ICES Journal of Marine Science*, fsx041. doi:
655 10.1093/icesjms/fsx041.

- 656 Paterson, S., Piertney, S. B., Knox, D., Gilbey, J., and Verspoor, E. 2004. Characterization and PCR
657 multiplexing of novel highly variable tetranucleotide Atlantic salmon (*Salmo salar* L.)
658 microsatellites. *Molecular Ecology Notes*, 4: 160-162.
- 659 Peakall, R. O. D., and Smouse, P. E. 2006. GENALEX 6: genetic analysis in Excel. Population genetic
660 software for teaching and research. *Molecular Ecology Notes*, 6: 288-295.
- 661 Pritchard, J. K., Stephens, M., and Donnelly, P. 2000. Inference of Population Structure Using
662 Multilocus Genotype Data. *Genetics*, 155: 945-959.
- 663 Rannala, B., and Mountain, J. L. 1997. Detecting immigration by using multilocus genotypes.
664 *Proceedings of the National Academy of Science, USA.*, 94: 9197-9201.
- 665 Rousset, F. 2008. genepop'007: a complete re-implementation of the genepop software for
666 Windows and Linux. *Molecular Ecology Resources*, 8: 103-106.
- 667 Saitou, N., and Nei, M. 1987. The neighbor-joining method: a new method for reconstructing
668 phylogenetic trees. *Molecular Biology and Evolution*, 4: 406-425.
- 669 Sanchez, J. A., Clabby, C., Ramos, D., Blanco, G., Flavin, F., Vazquez, E., and Powell 1996. Protein and
670 microsatellite single locus variability in *Salmo salar* L. (Atlantic salmon). *Heredity*, 77: 423-
671 432.
- 672 Seeb, L. W., Antonovich, A., Banks, M. A., Beacham, T. D., Bellinger, M. R., Blankenship, S. M.,
673 Campbell, M. R., et al. 2007. Development of a standardized DNA database for Chinook
674 salmon. *Fisheries*, 32: 540-552.
- 675 Shaklee, J. B., Beacham, T. D., Seeb, L., and White, B. A. 1999. Managing fisheries using genetic data:
676 case studies from four species of Pacific salmon. *Fisheries Research*, 43: 45-78.
- 677 Shedd, K. R., Dann, T. H., Hoyt, H. A., Foster, M. B., and Habicht, C. 2016. Genetic Baseline of North
678 American Sockeye Salmon for Mixed Stock Analyses of Kodiak Management Area
679 Commercial Fisheries, 2014–2016. Fishery Manuscript Series No. 16-03. Alaska Department
680 of Fish and Game. Anchorage, Alaska. 233 pp.
- 681 Sheehan, T. F., Legault, C. M., King, T. L., and Spidle, A. P. 2010. Probabilistic-based genetic
682 assignment model: assignments to subcontinent of origin of the West Greenland Atlantic
683 salmon harvest. *ICES Journal of Marine Science*, 67: 537-550.
- 684 Stephenson, J. J., Campbell, M. R., Hess, J. E., Kozfkay, C., Matala, A. P., McPhee, M. V., Moran, P., et
685 al. 2009. A centralized model for creating shared, standardized, microsatellite data that
686 simplifies inter-laboratory collaboration. *Conservation Genetics*, 10: 1145-1149.
- 687 Takezaki, N., Nei, M., and Tamura, K. 2010. POPTREE2: Software for Constructing Population Trees
688 from Allele Frequency Data and Computing Other Population Statistics with Windows
689 Interface. *Molecular Biology and Evolution*, 27: 747-752.
- 690 Tornaiainen, J., Vuorinen, P. J., Jones, R. I., Keinänen, M., Palm, S., Vuori, K. A. M., and Kiljunen, M.
691 2013. Migratory connectivity of two Baltic Sea salmon populations: retrospective analysis
692 using stable isotopes of scales. *ICES Journal of Marine Science*, 71: 336-344.
- 693 Vähä, J.-P., Erkinaro, J., Falkegård, M., Orell, P., and Niemelä, E. 2016. Genetic stock identification of
694 Atlantic salmon and its evaluation in a large population complex. *Canadian Journal of
695 Fisheries and Aquatic Sciences*: 1-12.
- 696 Verspoor, E., Beardmore, J. A., Consuegra, S., Garcia de Leaniz, C., Hindar, K., Jordan, W. C., Koljonen,
697 M. L., et al. 2005. Population structure in the Atlantic salmon: insights from 40 years of
698 research into genetic protein variation. *Journal of Fish Biology*, 67: 3-54.
- 699 Verspoor, E., Consuegra, S., Fridjonsson, O., Hjørleifsdóttir, S., Knox, D., Olafsson, K., Tompsett, S., et
700 al. 2012. Regional mtDNA SNP differentiation in European Atlantic salmon (*Salmo salar*): an
701 assessment of potential utility for determination of natal origin. *ICES Journal of Marine
702 Science*, 69: 1625-1636.
- 703 Waples, R. S., Kalinowski, S. T., and Anderson, E. C. 2008. An improved method for predicting the
704 accuracy of genetic stock identification. *Canadian Journal of Fisheries and Aquatic Sciences*,
705 65: 1475-1486.
- 706

This is the author's accepted manuscript. The final published version of this work (the version of record) is published by Oxford Academic in ICES Journal of Marine Science available at: <https://doi.org/10.1093/icesjms/fsx184>. This work is made available online in accordance with the publisher's policies. Please refer to any applicable terms of use of the publisher.

707

708

709

710 Table 1. Sample baseline coverage pre- and post-genotype quality control (see text for
711 details).

Country	Pre-QC			Post-QC		
	Rivers	Sites	Fish	Rivers	Sites	Fish
Denmark ¹	3	6	253	2	4	189
England ^{2,3}	24	38	1652	23	35	1498
Finland ⁴	2	5	395	2	5	393
France ^{2,3,5,6}	13	16	759	9	9	450
Iceland ⁷	17	25	2352	16	22	1986
Ireland ⁸	29	45	2345	29	40	2053
Northern Ireland ⁹	9	20	1469	7	18	1302
Norway ^{4,10,11}	90	109	7749	81	99	7008
Russia ^{4,10,12}	33	36	2506	30	33	2350
Scotland ³	87	230	11625	69	185	8884
Spain ⁶	7	7	342	4	4	190
Sweden ^{1,4}	4	4	180	4	4	172
Wales ²	7	10	375	6	9	347
Total	325	551	32002	282	467	26822

712 Institutions contributing data: ¹ Danish Institute for Fisheries Research, Denmark; ² University of Exeter,
713 England; ³ Marine Scotland Science, Scotland; ⁴ University of Turku, Finland; ⁵ Geneindex, France; ⁶ University
714 of Oviedo, Spain; ⁷ Marine and Freshwater Research Institute, Iceland; ⁸ University College Cork, Ireland; ⁹
715 Queen's University Belfast & Agri-Food and Biosciences Institute Northern Ireland, Northern Ireland; ¹⁰
716 Institute of Marine Research, Norway; ¹¹ Norwegian Institute for Nature Research, Norway, ¹² Knipovich Polar
717 Research Institute of Marine Fisheries & Oceanography, Russia.

718
719
720

721 Table 2. Individual assignment accuracy using fish removed from the reference baseline. Initial assignment units in first column defined by
 722 distance and STRUCTURE based analysis. Remaining assignments represent amalgamations of units where assignment accuracy is <80%.
 723 Assignment accuracy was calculated using only fish with individual assignment probabilities ≥ 0.8 . Values in bold represent accuracy of at least
 724 80% to assignment units. Sample sizes represent baseline/mixture size.

725

Assignment unit	Sample size	Assigned %	Correct %	Assignment unit	Assigned %	Correct %	Assignment unit	Assigned %	Correct %
White Sea	758/86	68.6	90.2	White Sea	70.9	90.3	White Sea	72.1	90.3
Kola	1561/160	50	82.1	Kola	51.9	82.1	Kola	53.1	82.1
Kola (Tuloma Basin)	287/39	61.5	100	Kola (Tuloma Basin)	66.7	96	Kola (Tuloma Basin)	66.7	96
Finnmark	1109/107	54.2	84.7	Finnmark	59.3	82.9	Finnmark	59.3	82.9
Teno/Tana	271/28	42.9	10						
Mid Norway	3195/369	54.5	84.1	Mid & SW Norway	68.3	84.4	Mid & SW Norway	69.2	84.4
South West Norway	816/95	42.1	73.8						
South Norway	693/83	32.5	81.25	South Norway	45.8	82.4	South Norway	47.0	82.4
Enningdalselva	86/8	87.5	100	Enningdalselva	87.5	100	Enningdalselva	87.5	100
Sweden	108/12	33.3	100	Sweden	41.7	100	Sweden	41.7	100
Baltic	47/5	60	100	Baltic	80.0	100	Baltic	80.0	100
Denmark	176/13	61.5	100	Denmark	76.9	100	Denmark	76.9	100
Central Scotland/North England	1711/200	32	73.5	Scotland/North East England	66.3	80.4	Scotland/NE England/Irish Sea	76.4	87.2
North East Scotland	2183/233	42.5	56.5						
Kyle/Ness	814/99	42.4	78.7						
North & West Scotland	2005/255	35.7	72						
Water of Luce	225/20	30	40						
West Central Scotland	242/28	46.4	83.3						
Irish Sea	1992/214	39.7	77.3	Irish Sea	52.3	76.3			
Leven	324/41	75.6	100	Leven	82.9	96.9	Leven	85.4	96.9
English Chalk	134/9	88.9	100	English Chalk	88.9	100	English Chalk	100	100
North France	283/35	45.7	78.9	North France	51.4	78.9	France/Spain	68.0	91.7
South France/Spain	282/40	70	100	South France/Spain	72.5	100			
North Ireland	1519/161	50.3	87.0	North Ireland	59.6	85.9	Ireland	64.0	87.0
South West Ireland	341/35	54.3	85.7	South West Ireland	52.2	77.4			
South Ireland	572/57	29.8	58.8						
Bann	619/51	66.7	93.9	Bann	66.7	93.9	Bann	66.7	93.9
North Iceland	976/110	95.5	96.3	North Iceland	95.5	96.3	North Iceland	95.5	96.3
West Iceland	811/89	91.0	98.7	West Iceland	92.1	98.7	West Iceland	93.3	98.7

726

727 Figure 1. Map of sampling region. Points represent sample sites and/or river mouths. Full
728 site information is contained in Supplementary data S1 and an expanded map with all rivers
729 identified is in Supplementary data S2. Regions noted are all those referred to in the text.
730 The Level one assignment units (see text) are delineated by the dashed line and the initial
731 Level 2 units by coloured points.

732

733 Figure 2. Neighbour-joining phylogenetic tree of sample sites based on D_A with major
734 clusters coloured and named. Expanded tree with all sites identified is detailed in
735 Supplementary data S4.

736

737 Figure 3. Hierarchical STRUCTURE based clustering analysis of sites. Each cluster analysis is
738 described using three components. Firstly the results of the STRUCTURE analysis are shown
739 with vertical bars representing individual sites and colours relating to cluster membership of
740 that site. A plot of the ΔK values (Earl and vonHoldt, 2012) associated with the analysis is
741 also shown defining the K identified in that cluster analyses. Finally a map is shown detailing
742 the geographic location of the clusters identified. Cluster names in italics refer to clusters for
743 which further hierarchical analysis was performed. Cluster names in regular text refer to
744 final cluster assignment groups.

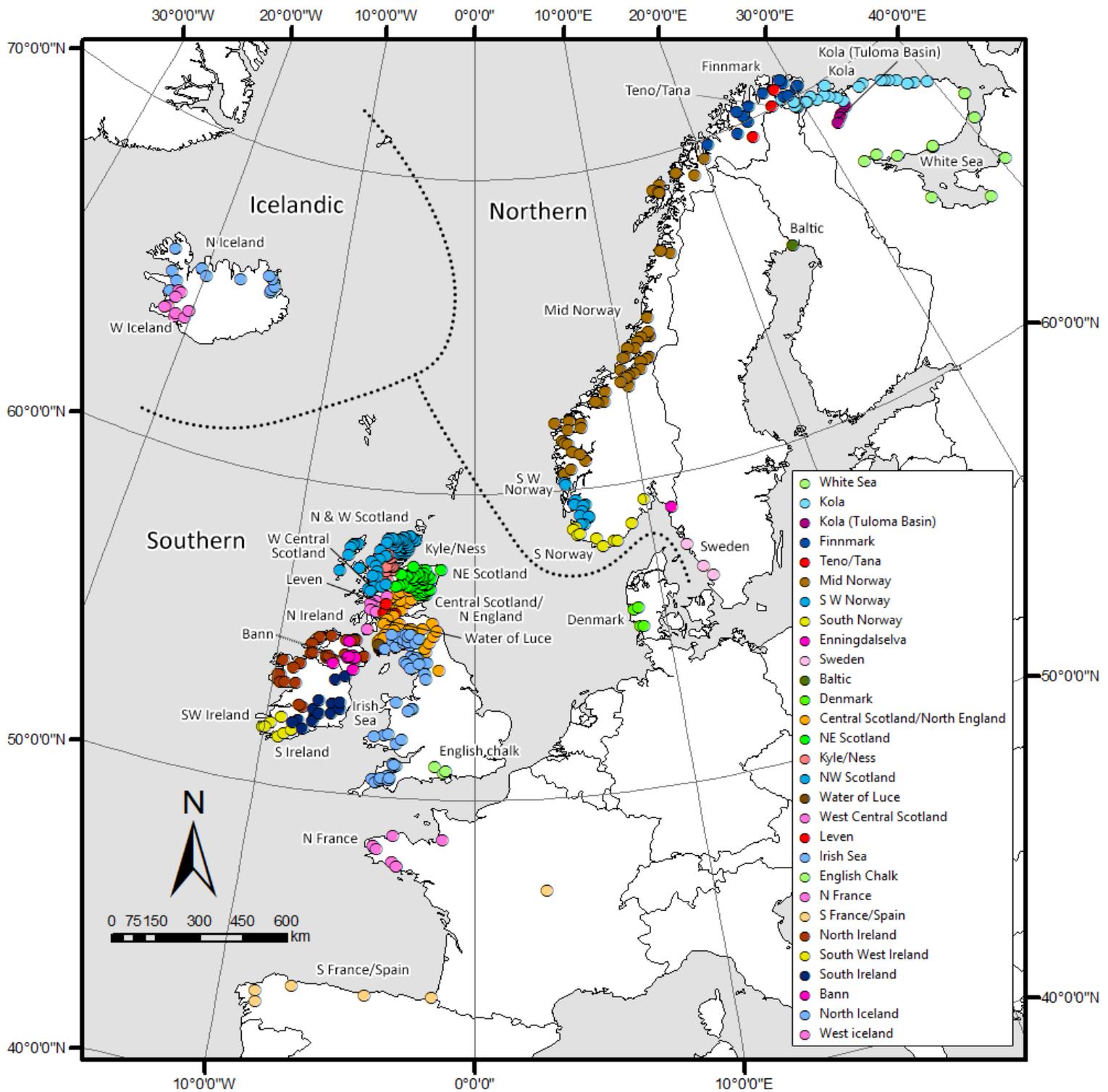
745

746 Figure 4. Proportion estimates from independent 100% simulation studies of the genetic
747 baseline at Level 1 and all stages of the iterative formation of the Level 2 assignment unit
748 levels. Points represent mean estimates with bars showing 95% confidence intervals.

749

750 Figure 5. A) Mixed stock fishery estimates using fish removed from the baseline and used as
751 fishery mixtures. B) Mixed stock fishery estimates using simulated equal proportions of fish
752 from each assignment unit in the mixture. C) Mixed stock fishery estimates using entire
753 rivers removed from the baseline and used as fishery mixtures. Dark bars represent actual
754 proportions in the mixture files and grey bars ONCOR estimates. Bars represent mean
755 estimates with 95% confidence intervals around these estimates. NOTE change of Y-axis
756 scale for the Level 1 and 2 assignment levels.
757

758 Fig. 1.



780 Fig. 2.

781

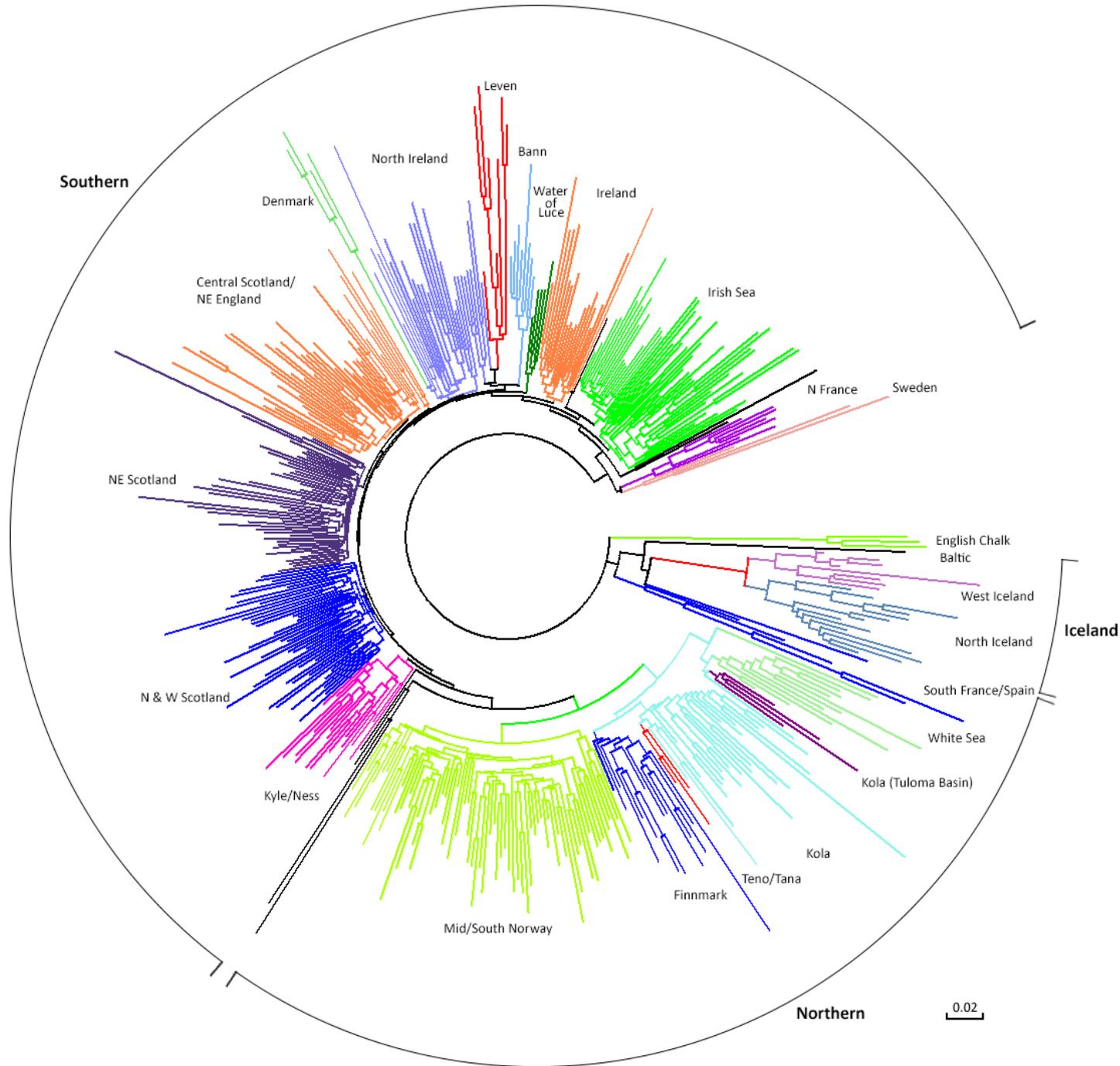
782

783

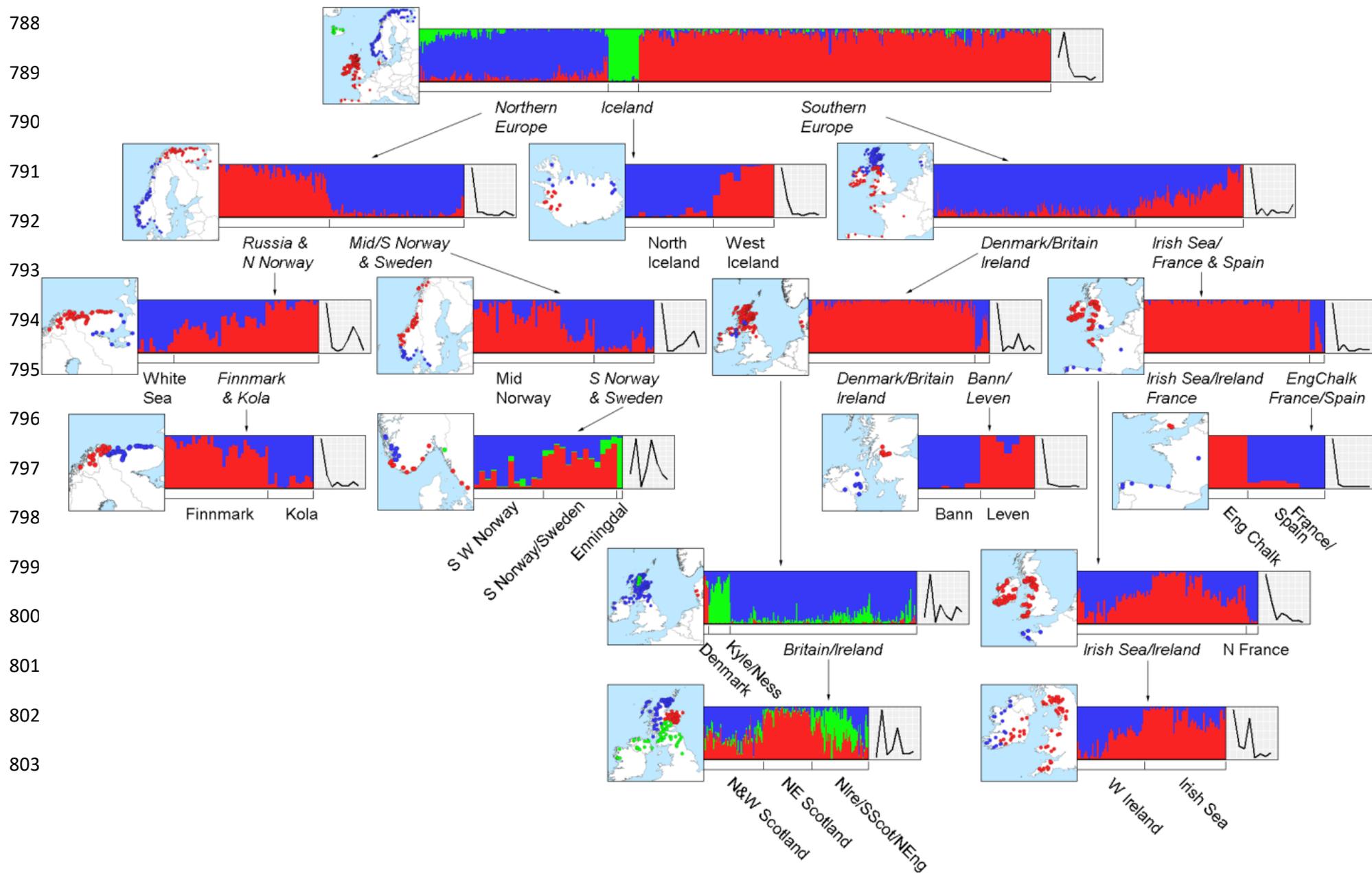
784

785

786

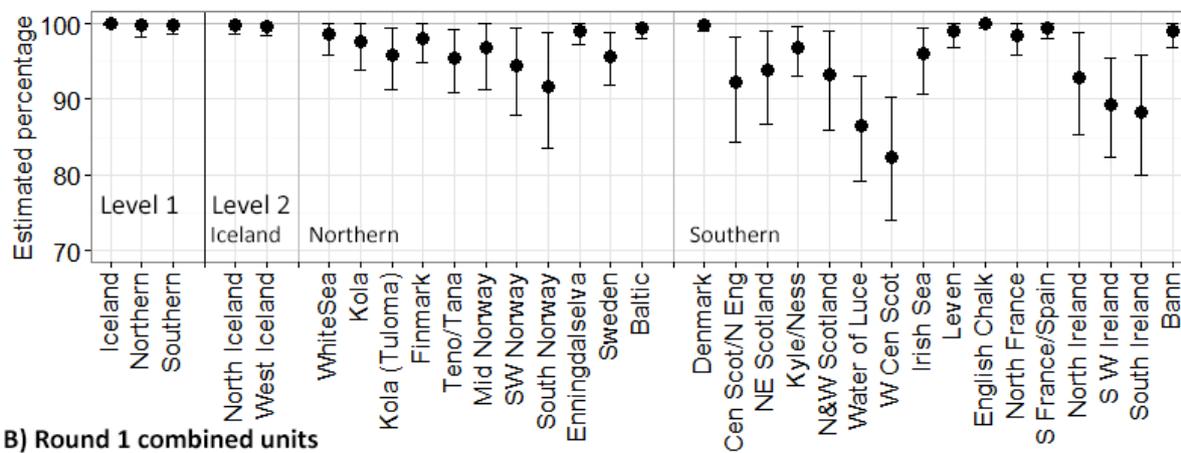


787 Fig. 3.

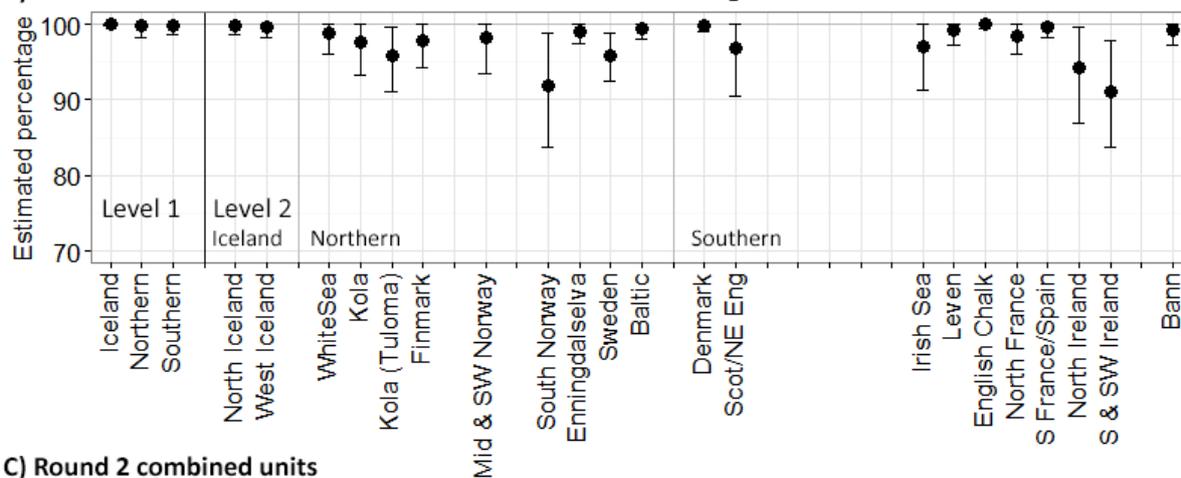


804 Fig. 4

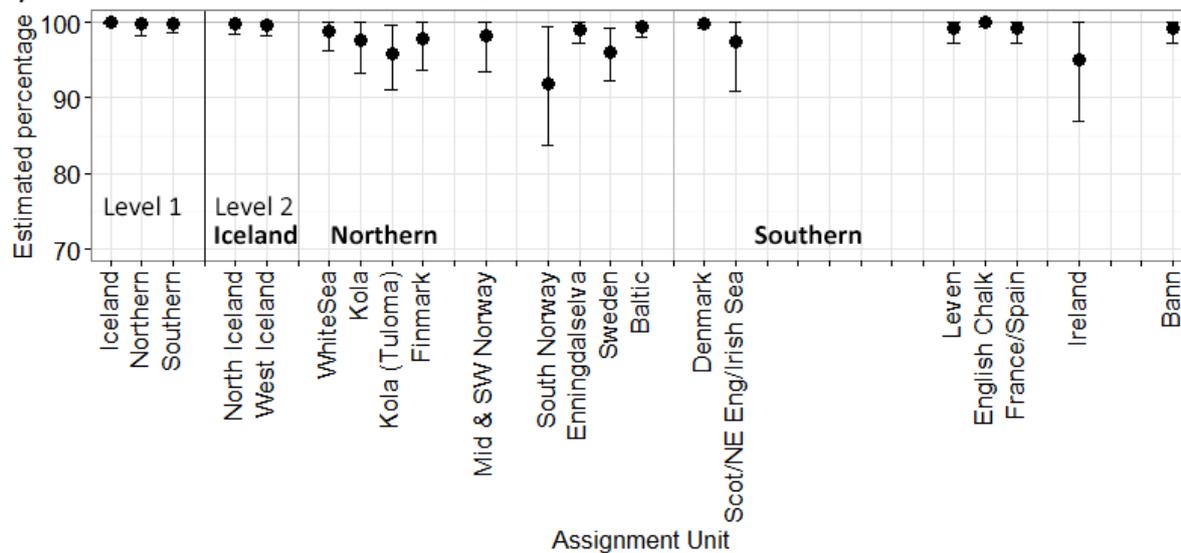
A) Initial assignment units



B) Round 1 combined units



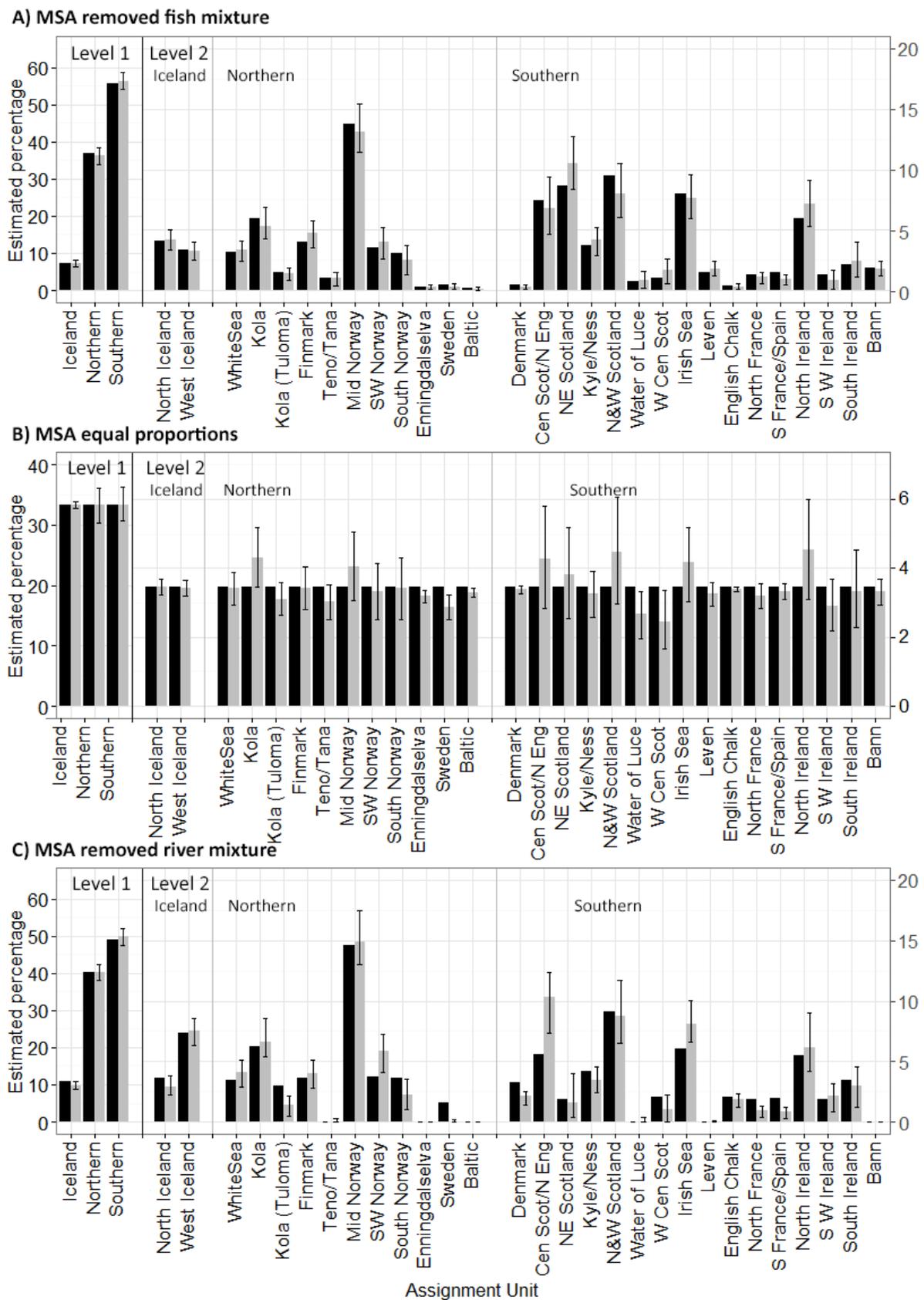
C) Round 2 combined units



805

806

807 Fig. 5



808
809