

# MOLECULAR ECOLOGY RESOURCES

## **In-field genetic stock identification of overwintering coho salmon in the Gulf of Alaska: Evaluation of Nanopore sequencing for remote real-time deployment**

Journal:	<i>Molecular Ecology Resources</i>
Manuscript ID	MER-21-0504.R1
Manuscript Type:	Resource Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Deeg, Christoph; The University of British Columbia, Department of Forest & Conservation Science; Pacific Salmon Foundation Sutherland, Ben; Fisheries and Oceans Canada, Pacific Biological Station Ming, Tobi; Fisheries and Oceans Canada, Pacific Biological Station Wallace, Colin; Fisheries and Oceans Canada, Pacific Biological Station Jonsen, Kim; Fisheries and Oceans Canada, Pacific Biological Station Flynn, Kelsey; Fisheries and Oceans Canada, Pacific Biological Station Rondeau, Eric; Fisheries and Oceans Canada - Pacific Biological Station Beacham, Terry; Fisheries and Oceans Canada, Pacific Biological Station Miller, Kristina; Fisheries and Oceans Canada, Pacific Biological Station; The University of British Columbia, Department of Forest & Conservation Science
Keywords:	Bioinformatics/Phyloinformatics, Conservation Biology, Fisheries Management, Nanopore Sequencing, Pacific Salmon, Mobile Genetic Stock Identification

# 1 **In-field genetic stock identification of overwintering coho** 2 **salmon in the Gulf of Alaska: Evaluation of Nanopore** 3 **sequencing for remote real-time deployment**

4  
5 Christoph M. Deeg<sup>1,2</sup>, Ben J. G. Sutherland<sup>3</sup>, Tobi J. Ming<sup>3</sup>, Colin Wallace<sup>3</sup>, Kim Jonsen<sup>3</sup>, Kelsey L. Flynn<sup>3</sup>, Eric B.  
6 Rondeau<sup>3</sup>, Terry D. Beacham<sup>3</sup>, and Kristina M. Miller<sup>1,3</sup>

7 1: Forest and Conservation Sciences, University of British Columbia, Vancouver, BC, Canada

8 2: Pacific Salmon Foundation, Vancouver, BC, Canada

9 3: Fisheries and Oceans Canada, Pacific Biological Station, Nanaimo, BC, Canada

## 10 **Abstract:**

11 Genetic stock identification (GSI) from genotyping-by-sequencing of SNP loci has become the gold  
12 standard for stock of origin identification in Pacific salmon. Sequencing platforms currently applied require large  
13 batch sizes and multi-day processing in specialized facilities to perform genotyping by the thousands. However,  
14 recent advances in third-generation single-molecule sequencing platforms, like the Oxford Nanopore minION,  
15 provide base calling on portable, pocket-sized sequencers and promise real-time, in-field stock identification of  
16 variable batch sizes. Here we evaluate utility and comparability to established GSI platforms of at-sea stock  
17 identification of coho salmon (*Oncorhynchus kisutch*) using targeted SNP amplicon sequencing on the minION  
18 platform during a high-sea winter expedition to the Gulf of Alaska. As long read sequencers are not optimized for  
19 short amplicons, we concatenate amplicons to increase coverage and throughput. Nanopore sequencing at-sea  
20 yielded data sufficient for stock assignment for 50 out of 80 individuals. Nanopore-based SNP calls agreed with Ion  
21 Torrent based genotypes in 83.25%, but assignment of individuals to stock of origin only agreed in 61.5% of  
22 individuals highlighting inherent challenges of Nanopore sequencing, such as resolution of homopolymer tracts and  
23 indels. However, poor representation of assayed salmon in the queried baseline dataset contributed to poor  
24 assignment confidence on both platforms. Future improvements will focus on lowering turnaround time and cost,  
25 increasing accuracy and throughput, as well as augmentation of the existing baselines. If successfully implemented,

26 Nanopore sequencing will provide an alternative method to the large-scale laboratory approach by providing mobile  
27 small batch genotyping to diverse stakeholders.

## 28 Key words:

29 Nanopore, genetic stock identification, single nucleotide polymorphism, salmon, at-sea, mobile

## 30 Introduction:

31 Pacific salmon are crucial to coastal and terrestrial ecosystems around the North Pacific by connecting  
32 oceanic and terrestrial food webs and nutrient cycles (Cederholm et al. 1999). Salmon are highly valued by the  
33 northern Pacific Rim nations due to their contribution to commercial and recreational fisheries as well as their  
34 cultural importance, especially amongst Indigenous peoples (Lichatowich 2001). Despite this significance, many  
35 wild Pacific salmon stocks have experienced population declines due to a combination of compounding factors such  
36 as overexploitation, spawning habitat alterations, pathogens and predators, prey availability, and climate change  
37 (Miller et al. 2014). Efforts to rebuild stocks include habitat restoration, artificial stock enhancements, as well as  
38 stock specific monitoring through several assessment methods to inform targeted management and harvest strategies  
39 (Hinch et al. 2012). Stock specific management can be implemented through traditional small scale terminal  
40 fisheries, but the majority fisheries occur in mixed stock environments where stock identification methods are crucial  
41 to minimize impact on stocks of concern while allowing the harvest of abundant stocks (Atlas et al. 2021; Dann et al.  
42 2013).

43 To inform mixed-stock management, stock identification has in the distant past utilized characteristic scale  
44 and parasite patterns as well as the marking of hatchery-enhanced fish by coded-wire tagging (Wood, Rutherford,  
45 and McKinnell 1989; Cook and Guthrie 1987; Jefferts, Bergman, and Fiscus 1963). More recently, genetic stock  
46 identification (GSI) using allozyme, minisatellite, microsatellite, and ultimately single nucleotide polymorphisms  
47 (SNPs) as markers has proven superior in delivering high-throughput insights into the stock composition of salmon  
48 (Winans et al. 1994; Miller, Withler, and Beacham 1996; Beacham et al. 2017, 2018). Specifically, the large baseline  
49 of population-specific SNP frequencies and targeted amplification of such SNP loci now allow for unprecedented  
50 resolution of stock origin in many species of salmon at reduced biases (Beacham et al. 2017, 2018; Ozerov et al.

51 2013; Gilbey et al. 2017). However, current sequencing approaches, based on second generation sequencing  
52 platforms (e.g. Illumina and Ion Torrent), mean that only sequencing large batches of individuals, known as  
53 “genotyping by the thousands” (GT-seq), is economically sensible (Beacham et al. 2017, 2018; Campbell, Harmon,  
54 and Narum 2015). These approaches require a specialized laboratory and several days turnover for the library  
55 preparation and sequencing, even under highly automated settings. These constraints limit the utility of SNP-based  
56 GSI for real world scenarios that are often spatially or temporally restricted, because samples need to be transported  
57 to the laboratory for analysis, as has been the case for most GSI methods to date. Specifically, for time-sensitive  
58 stock-specific harvest management decisions, an in-field real-time SNP-based GSI approach with greater flexibility  
59 in sample batch size would be desirable.

60         Recent advances in third-generation single-molecule sequencing platforms like the Oxford Nanopore  
61 minION allow real-time sequencing on a pocket-sized portable sequencer that requires little library preparation,  
62 therefore enabling sequencing in remote locations (Mikheyev and Tin 2014; Quick et al. 2016). However, several  
63 technical hurdles to adapting Nanopore sequencing to SNP GSI exist. While Nanopore sequencing can yield  
64 extremely long reads, the number of sequencing pores and their loading rate is limited, resulting in low throughput  
65 when sequencing short reads such as amplicons. An additional problem is the relatively high error rate inherent to  
66 this novel technology. Since the SNP GSI protocols are based on the amplification of short amplicons via targeted  
67 multiplex PCR, sequencing throughput of such short amplicons on the Nanopore platform is comparatively low, as  
68 the number of sequencing pores is the rate limiting factor. This is especially problematic since high coverage is  
69 needed to compensate for the higher error rate of Nanopore generated sequences. A promising approach to overcome  
70 these limitations is the concatenation of PCR amplicons that allows the sequencing of several amplicons within a  
71 single read, thereby exponentially increasing throughput for genotyping (Cornelis et al. 2017; Schlecht et al. 2017).

72         Here, we report on the development and performance of a novel Nanopore-based in-field SNP GSI method  
73 by adapting existing SNP GSI technology to the Nanopore platform using a concatenation approach (Schlecht et al.  
74 2017). We aim to demonstrate in-field feasibility, repeatability, and comparability to established platforms. As a  
75 proof of concept, in-field stock ID was performed in the Gulf of Alaska onboard the research vessel *Professor*  
76 *Kaganovsky* during the International Year of the Salmon (IYS) expedition in February and March of 2019.

## 77 Materials and Methods:

### 78 Field Lab equipment and workspace

79 The field equipment onboard the *Professor Kaganovsky* research trawler consisted of a PCR thermocycler, a  
80 mini-plate centrifuge, a microcentrifuge, a Qubit fluorimeter (Thermo Fisher), a vortexer, a minION sequencer, a  
81 laptop with an Ubuntu operating system (Ubuntu v.14.06), as well as assorted pipettes and associated consumables  
82 like filter tips (Figure 1). The required infrastructure onboard included a 4°C fridge, a -20°C freezer, power supply,  
83 as well as a physical workspace. The entire equipment configuration required was under \$10,000 CAD.

### 84 Tissue sample collection and DNA extraction

85 Salmon were captured by the research trawler *Professor Kaganovsky* during the 2019 International Year of  
86 the Salmon (IYS) Signature expedition in the Gulf of Alaska (Supplemental Figure 1). We collected fin clips of coho  
87 salmon (*Oncorhynchus kisutch*) and froze them individually until DNA extraction, or immediately processed once a  
88 suitable batch size had been accumulated. DNA extraction from 2 x 2 x 2mm fin-tissue clips was performed in a 96-  
89 well PCR plate using 100µl of QuickExtract solution (Lucigen, USA) according to the manufacturer's instructions.

### 90 Multiplex PCR and Barcoding

91 Multiplex PCR with a custom panel of primers targeting 299 loci of known SNPs was performed using  
92 0.25µl of DNA extract as template using the AgriSeq HTS Library Kit Amplification Mix PCR mastermix  
93 (ThermoFisher) in a 10µl reaction according to Beacham et al. (Beacham et al. 2017; See appendix A2). Primer sets  
94 targeting Multi nucleotide polymorphisms (MNPs) were included in the primer panel by Beacham et al. 2017, but  
95 were excluded from the analysis (Supp. Table 1). Next, we prepared amplicons for ligation by end-prepping  
96 amplified strands with AgriSeq HTS Library Kit Pre-ligation Enzyme. ONT barcode adapters (PCR Barcoding  
97 Expansion 1-96, EXP-PBC096, Oxford Nanopore Technologies, UK) were then ligated to the amplicons by blunt-  
98 end ligation with the Barcoding Enzyme/Buffer of the AgriSeq HTS Library Kit according to manufacturer's  
99 instruction. After bead-cleanup (1.2:1 bead:sample, AMPure XP beads, Beckman Coulter, USA) we added the  
100 ligation products, barcodes and barcoding adapters (PCR Barcoding Expansion 1-96, EXP-PBC096, Oxford  
101 Nanopore Technologies, UK) by PCR using Q5 polymerase mastermix (NEB, USA) for individual fish identification

102 according to manufacturer's protocol in a 25µl reaction (98°C for 3 min; 25 cycles of 98°C for 10 s, 70°C for 10 s,  
103 72°C for 25 s; 72°C for 2 min). Barcoded libraries were then pooled and cleaned using 1.2:1 bead cleanup, before  
104 DNA yield of a subset of samples (12.5%) was analyzed by Qubit (dsDNA HS Assay Kit, ThermoFisher, USA).

## 105 Amplicon concatenation

106 To improve throughput on the minION, we concatenated amplicons using inverse complementary adapters  
107 (Figure 2). After end prep using Ultra II End Repair/dA-Tailing Module Module (NEB, USA), the library was split  
108 into two equal volume subsets. Custom inverse complementary adapters that had inverse complementary terminal  
109 modifications to ensure unidirectional ligation (3'-T overhang and 5' phosphorylation) were ligated onto both ends of  
110 the respective subsets using the Ultra II Ligation Module (NEB, USA) according to manufacturer's instructions and  
111 purified with 1:1 bead cleanup (Figure 2). The custom adapters were adapted from Schlecht et al (Schlecht et al.  
112 2017): Adapter A: 5'P-ACAGCGAGTTATCTACAGGTTCTTCAATGT +  
113 ACATTGAAGAACCTGTAGATAACTCGCTGTT ;  
114 Adapter B: 5'P-ACATTGAAGAACCTGTAGATAACTCGCTGT +  
115 ACAGCGAGTTATCTACAGGTTCTTCAATGTT). Amplicons with adapters added to them were subsequently  
116 amplified again with a single primer (ACATTGAAGAACCTGTAGATAACTCGCTGTT for adapter A,  
117 ACAGCGAGTTATCTACAGGTTCTTCAATGTT for adapter B) in 25µl Q5 reactions according to manufacturer's  
118 instructions with the following thermal regime: 98°C for 3 min; 30 cycles of 98°C for 10 s, 68°C for 15 s, 72°C for  
119 20 s; 72°C for 2 min). After 1:1 bead cleanup, we pooled both subsets in equimolar ratios after Qubit quantification  
120 to verify both reactions worked, and then subjected the pool to a primer-free, PCR-like concatenation due to  
121 heterodimer annealing and elongation in 25µl Q5 reaction, using the complementary adapter sequence ligated onto  
122 the amplicons as primers cycled under the following thermal regime: 3 cycles of 98°C for 10 s, 68°C for 30 s, 72°C  
123 for 20 s; followed by 3 cycles of 98°C for 10 s, 68°C for 30 s, 72°C for 30 s; followed by 3 cycles of 98°C for 10 s,  
124 68°C for 30s,72°C for 40s; followed by 3 cycles of 98°C for 10 s, 68°C for 30 s, 72°C for 50s; and finally followed  
125 by 72°C for 2 min (Figure 2).

## 126 Library Preparation and sequencing

127           The concatenated amplicons were prepared for Nanopore sequencing using the ONT Ligation Sequencing  
128 Kit (LSK109) according to the manufacturer's instruction. In brief, after end-prep using the Ultra II Endprep Module  
129 and bead cleanup, we ligated proprietary ONT sequencing adapters onto the concatenation adapters by blunt-end  
130 ligation using the proprietary ONT Buffer and the TA quick ligase (NEB, USA; note: this standard sequencing step  
131 not shown in Figure 2). After additional bead-cleanup and washing with the short fragment buffer (SFB: ONT, UK)  
132 according to the manufacturer's protocol, we loaded the library onto a freshly primed flow cell (MIN 106 R9.4.1:  
133 ONT, UK) according to the manufacturer's instruction.

## 134 Nanopore sequencing, deconcatenation, and binning

135           After flow cell priming and loading of the library, the flow cell was placed on the minION sequencer.  
136 Sequencing and basecalling into fast5 and fastq was performed simultaneously using minKNOW (version 3.1.8) on  
137 an Ubuntu 14.06 platform. First, all fastq raw reads that passed default quality control in minKNOW were combined  
138 into bins of 500k reads each. This had empirically been determined to be the maximum number of reads allowing  
139 simultaneous processing in the downstream analysis on our platform (Ubuntu 14.06, 31.2 GiB RAM 7700K CPU @  
140 4.20GHz × 8). Reads containing concatenated amplicons were deconcatenated and the concatenation adapter  
141 sequence was trimmed off the remaining sequence using porechop (<https://github.com/rrwick/Porechop>) with a  
142 custom adapter file (“adapters.py”) that only contained the concatenation adapter under the following settings:  
143 `porechop-runner.py -i input_raw_reads.fastq -o output/dir -t 16 --middle_threshold 75 --min_split_read_size 100 --`  
144 `extra_middle_trim_bad_side 0 --extra_middle_trim_good_side 0`

145 We binned the deconcatenated reads by barcode corresponding to fish individuals by using porechop with the  
146 provided default adapters file and the following settings:

147 `porechop-runner.py -i input_deconcatenated_reads.fastq -b binning/dir -t 16 --adapter_threshold 90 --end_threshold`  
148 `75 --check_reads 100000`

149 After this step, all reads from the corresponding barcode bins corresponding to the same individual across the  
150 different 500k sub-bins were combined for downstream analysis. See <https://github.com/bensutherland/nano2geno/>  
151 for source scripts for analysis.

## 152 Alignment and SNP calling

153 We aligned the binned reads to the reference amplicon sequences described by Beacham et al. 2017 using  
154 BWA-MEM and indexed using samtools (Beacham et al. 2017; Li et al. 2009; Li and Durbin 2009). Alignment  
155 statistics for all loci were generate using pysamstats (<https://github.com/alimanfoo/pysamstats>; flags: -t variation -f)  
156 and we extracted the nucleotides observed at the relevant SNP hotspot loci from the resulting file using a custom R  
157 script by looping through the results file guided by a SNP location file. Finally, we compared the observed  
158 nucleotide distributions at SNP hotspots with to the hotspot reference and variant nucleotides and scored as  
159 homozygous reference when  $\geq 66\%$  of the nucleotides were the reference allele, heterozygous when the reference  
160 allele was present  $< 66\%$  and the variant allele  $> 33\%$ , or as homozygous variant (when the nucleotides were  $\geq 66\%$   
161 the variant allele) using a custom R script to generate a numerical locus table. We visually inspected alignments  
162 determined to be problematic using the IGV viewer (Robinson et al. 2011). The full pipeline titled “nano2geno”  
163 (n2g) including all custom scripts can be found at <https://github.com/bensutherland/nano2geno/> (Figure 2).

## 164 Mixed-stock Analysis

165 We performed mixture compositions and individual assignments using the R package rubias (Moran and  
166 Anderson 2019) with default parameters against the coho coastwide baseline of known allele frequencies for these  
167 markers established by Beacham et al (Beacham et al. 2017, 2020). The baseline used in this manuscript is available  
168 at <https://doi.org/10.5061/dryad.g4f4qrfs3>.

## 169 Ion torrent sequencing

170 To confirm the results obtained by Nanopore sequencing, the samples were sequenced using an Ion Torrent  
171 sequencer according to Beacham et al. 2017. In brief: DNA was extracted from the frozen tissue samples using the  
172 Biosprint 96 SRC Tissue extraction kit, and multiplex PCR and barcoding with Ion Torrent Ion Codes was  
173 performed using the AgriSeq HTS Library Kit (ThermoFisher). The libraries were then prepared with the Ion Chef  
174 for sequencing on the Ion Torrent Proton Sequencer and SNP variants were either called by the Proton VariantCaller  
175 (ThermoFisher; Torrent Suite 5.14.0) software or the custom SNP calling script of the nano2geno pipeline. The  
176 resulting locus score table was then analyzed using rubias as described above.

## 177 Concordance assessment

178 We assessed concordance between sequencing platforms on SNP level. A PCoA analysis was performed  
179 using the R package ape based on a reference vs allele call matrix using a restricted dataset including only  
180 individuals that had stock assignment on both platforms (Paradis and Schliep 2019). Additionally, calls (reference vs.  
181 alternate allele) were compared for each sample and marker individually, then averaged by individual, and then  
182 averaged by the entire assessed population. Similarly, we compared stock assignment by rubias by comparing the  
183 reporting unit or collection as assigned and scoring a match (1) or non-match (0). These scores were then averaged  
184 again to generate the final concordance or repeatability score as a percentage.

## 185 Results

### 186 In-field Nanopore Sequencing:

187 During the International Year of the Salmon Signature expedition to the Gulf of Alaska in February and  
188 March 2019, in-field single nucleotide polymorphism genetic stock identification (SNP GSI) was performed on coho  
189 salmon as the tissues became available. A total of 75 coho salmon were analyzed in two sequencing runs at different  
190 points during the expedition, representing 77% of all coho salmon captured during the expedition.

191 The first sequencing run was performed on February 26<sup>th</sup> and included 31 individuals. Library preparation onboard  
192 the vessel took 14h. However, faulty flow cell priming resulted in only approximately half the detected pores being  
193 active (843 pores). Of these pores, no more than 25% were actively sequencing at any time, highlighting the  
194 challenges of utilizing sensitive equipment under field conditions including excessive ship movement. Accordingly,  
195 sequencing for 30h and base-calling for 34h resulted in only 1.44M reads, 49% of which passed quality control. The  
196 read length distribution showed several large, concatenated amplicons up to 7,095 bp with a mean length of 825 bp  
197 (Supplemental Figure 2). Deconcatenation resulted in a read inflation by a factor of 2x (702k to 1,444k reads). After  
198 binning, reads per individual ranged from 1,983 to 86,467 reads with a mean of 13,709 reads (SD: 15,370), and  
199 722,174 reads that were not able to be assigned (50% of total deconcatenated reads) (Figure 3, Supplemental Figure  
200 2, Supplemental Figure 3).

201           The second sequencing run was performed on March 10<sup>th</sup>, 2019, with 44 coho salmon. Library preparation  
202 again took 14h and sequencing on a new flow cell took 15h, starting with 1,502 available pores, and up to 65%  
203 actively sequencing pores, and resulted in 4.48M reads, 76% of which passed quality control. Read lengths averaged  
204 810 bp with a maximum length of 8,023 bp (Supplemental Figure 2). Due to the large number of reads and the  
205 limited power of the computer being used for the analysis, base-calling into fastq took three days. Deconcatenation  
206 resulted in a read inflation of a factor of 1.7x (3.4M to 5.8M) (Supplemental Figure 2). Reads per individual showed  
207 a mean of 67,636 reads (SD: 59,393; min: 11,684; max: 335,348), with 722,179 reads remaining unassigned (12%)  
208 (Figure 3, Supplemental Figure 2, Supplemental Figure 3).

209           Upon return from the expedition, we sequenced 80 individuals, including all those previously genotyped  
210 aboard the vessel, in a single MinION run using the expedition setup starting from the frozen tissues from the  
211 expedition. We sequenced for 42h to maximize the total number of reads with 60% of 2,048 available pores actively  
212 sequencing resulting in 5.32 M reads. Of these reads, 3.20 M passed quality control. Again, large, concatenated  
213 amplicons up to 9,449 kb were observed, with a mean read length of 840 bp, and deconcatenation resulted in 4.54 M  
214 reads (1.4x inflation) (Supplemental Figure 2). The mean number of reads per bin was 29,439 (SD: 25,000) and  
215 ranged from 2,969 to 128,718 reads per individual, with 1,413,626 unassigned reads (31%) (Figure 3, Supplemental  
216 Figure 2, Supplemental Figure ).

217           Despite the absence of normalization between samples prior to multiplex PCR, barcoding, and loading, the  
218 binning distribution across samples was relatively even with only a few apparent outliers observed (Figure 3,  
219 Supplemental Figure 3). The minimum number of reads per individual sample necessary to cover sufficient loci (at a  
220 minimum depth of 10 sequences per locus) for downstream stock assignments (i.e., at least 141 loci per sample) is  
221 around 2,000 reads (Figure 3, Supplemental Figure 3).

## 222 Nanopore sequencing data requires loci reassessment for efficient SNP calling

223           After alignment to the reference sequences for SNP calling, Nanopore sequence data showed a  
224 comparatively higher error rate than Ion Torrent reads, as expected, with abundant indels that frequently led to lower  
225 alignment scores than those obtained by the Ion Torrent data (Ion Torrent average alignment score: 25.6 MAPQ;  
226 Nanopore average alignment score: 13.9 MAPQ). Specifically, regions containing homopolymer tracts were poorly  
227 resolved, as had previously been reported (Cornelis et al. 2017). Several instances could be identified where the

228 homopolymer presence near the SNP locus caused problematic alignments and therefore resulted in SNP calls not  
229 matching those found by the Ion Torrent on the same individual (Figure 4). Accordingly, six such loci were excluded  
230 from downstream analysis (Supp. Table 1). Other loci were excluded from the analysis due to absence of coverage  
231 (four loci) or the inability of the custom n2g pipeline to call MNPs (multi-nucleotide polymorphisms) or deletions  
232 (seven loci), bringing the number of accessed loci from 299 to 282 loci. Other loci showing apparent differences  
233 between Nanopore and Ion Torrent sequence data ( $n = 21$ ) were retained as no apparent explanation for the  
234 discrepancies could be identified.

235 After the removal of the discrepancies due to MNP, homopolymer, or deletion presence, the SNP cutoff for  
236 downstream analysis was set to 141 loci (50%). Only nine of 31 individuals (29%) of the first IYS sequencing run  
237 with problematic flow cell priming passed this threshold. In the second IYS sequencing run, 43 of 44 individuals  
238 passed the threshold (98%). The repeat run performed at the Pacific Biological Station resulted in 50 of the 80 (63%)  
239 that passed this threshold (Figure 3).

## 240 Platform biases lead to moderately altered SNP calling compared to Ion Torrent 241 sequencing

242 To assess the discrepancies between sequencing platforms, individuals that passed the genotyping rate  
243 threshold of 141 called loci (50% genotyping rate) in all data sets (i.e., Nanopore data during the expedition analyzed  
244 with n2g: “nano IYS”, Nanopore acquired during the repeat run upon return from the expedition, analyzed with n2g:  
245 “nano PBS”, Ion Torrent sequencing data analyzed with variant caller: “ion vc”, Ion Torrent analyzed with n2g: “ion  
246 n2g”) were included in a PCoA analysis on the SNP genotypes (Figure 5). This comparison excluded the MNP,  
247 deletion, and homopolymer loci (see above), but retained those without an explanation as to why the genotyping did  
248 not match. However, there was still an apparent separation by sequencing platform across the highest-scoring  
249 dimension (Figure 5). This platform dependent difference was reflected by 83.9% of SNP calls generated by  
250 Nanopore sequencing during the IYS expedition (nano IYS) and 83.7% of SNP calls generated during the repeat run  
251 upon return (nano PBS) matching the SNP calls based on Ion Torrent data (ion n2g) with nanopore reads having  
252 higher proportion of heterozygotes compared to Ion Torrent data (43% vs 33%). The agreement on SNP call between  
253 both Nanopore runs (comparing reference or alternate scores for both alleles from nano IYS vs nano PBS) was

254 84.4%, highlighting the inter run variability associated with current Nanopore sequencing. There was a slight  
255 correlation observed between the number of Nanopore reads per individual and the concordance with Ion Torrent  
256 SNP calls, suggesting that read depth is only a minor factor influencing SNP call concordance at the current  
257 threshold of a minimal alignment depth of 10x per site for Nanopore reads (Supplemental Figure 4). Excluding  
258 MNPs, deletions, and homopolymer issues, the influence of the SNP calling pipeline (n2g vs. variant caller) appears  
259 negligible compared to the differences by sequencing platform (Figure 5). Accordingly, SNPs scored based on the  
260 same Ion Torrent data sequence matched in 99.21% of cases between the two genotyping pipelines.

261 **Stock assignment based on Nanopore data is moderately repeatable and differs inherently**  
262 **from Ion Torrent based assignments in a subset of individuals**

263 Stock assignment by rubias showed discrepancies between the Nanopore and Ion Torrent based datasets. In  
264 only 61.5% of cases did Nanopore sequences (PBS run) lead to the same top reporting unit (reunit; large scale  
265 geographic areas such as Westcoast Vancouver Island or Lower Fraser River) assignment for individual stock ID as  
266 the Ion Torrent based sequences (Figure 6, Table 1). Specifically, Nanopore-based reunit assignment showed higher  
267 proportions of assignments to Southeastern Alaska (SEAK) than Ion Torrent-based assignments (Figure 6, Table 1).  
268 Nevertheless, mixture proportions in both datasets were dominated by Southeastern Alaska stocks. Nanopore  
269 assignments tended to overestimate the contribution to this stock as well as Lower Stikine River stocks (LSTK).  
270 Many of the individuals assigned to these stocks using the Nanopore were assigned to the adjacent stocks of Lower  
271 Hecate Strait and Haro Strait (HecLow+HStr) as well as Southern Coastal Streams, Queen Charlotte Strait, Johnston  
272 Strait and Southern Fjords (SC + SFj) on the Ion Torrent platform (Figure 6, Table 1). Individuals from stocks well  
273 represented in the database like the Columbia River were confidently assigned to the appropriate stock on both  
274 platforms. However, Z-scores calculated by rubias during stock assignment, which are an indirect measure of how  
275 well the SNP call match individuals in the baseline dataset of both, indicated that the Nanopore and the Ion Torrent  
276 data showed large deviations from the normal distribution, suggesting that many individuals assayed are not well  
277 represented in the database (Supplemental Figure 5) (Moran and Anderson 2019). Ion Torrent data shows two peaks,  
278 one overlaying the expected normal distribution and a second peak that lay outside of the normal distribution. This  
279 suggests that about half of the individuals were not from populations that are well represented in the database

280 (Supplemental Figure 5). Similarly, Nanopore-based assignments showed even more aberrant distribution,  
281 presumably due to the additive effects of the sequencing platform introducing bias on top of poor baseline  
282 representation (Supplemental Figure 5). The poor database representation could cause small differences in SNP calls  
283 to cause alternative assignments.

## 284 Discussion

285 Nanopore sequencing enables remote in-field single nucleotide polymorphism genetic  
286 stock identification

287 Here, we present the first proof-of-concept study demonstrating the feasibility of using the portable Oxford  
288 Nanopore minION sequencer for remote in-field genetic stock identification by SNP sequencing of Pacific salmon.  
289 We developed a rapid sample processing workflow that relied on amplicon concatenation to increase throughput.  
290 With this workflow, we performed genetic stock identification on 75 coho salmon onboard a research vessel in the  
291 Gulf of Alaska, with minimal equipment during two runs. Genetic stock identification of all 80 captured coho salmon  
292 in a single run using the mobile platform resulted in stock assignment for 50 individuals at 67% concordance with  
293 state of the art laboratory based pipelines.

294 Despite its promising performance, the fidelity, throughput, and turnaround time of Nanopore-based SNP  
295 GSI currently still falls short of what would enable this technology to be used for the wide range of remote real-time  
296 applications we intended it for. This is due to a number of factors, such as inefficient barcoding, error rates,  
297 inefficiencies of custom genotyping pipelines, low concatenation efficiency, and limited computational power in our  
298 setup. Further, the present protocol requires a high level of molecular laboratory expertise to perform the analysis.

299 The inherent low fidelity of the Nanopore platform using R9 type flow cells relative to other sequencing  
300 technologies, specifically around homopolymer tracts, proved to be the major shortcoming, limiting both the actual  
301 SNP calling accuracy, causing comparatively low repeatability, as well as the throughput, by necessitating a higher  
302 alignment coverage due to the high error rate (Cornelis et al. 2017). The low fidelity of the Nanopore sequences was  
303 specifically apparent when comparing it with the established sequencing platform for genetic stock identification by  
304 SNP sequencing, the Ion Torrent Proton sequencer (Beacham et al. 2017). The Ion Torrent short read sequencer

305 routinely outperformed the Nanopore sequencer, both in accuracy and in throughput. The latter being a major  
306 restricting factor of the Nanopore platform due to a limited number of available sequencing pores inherent to the  
307 platform. While we compensated for this limitation by concatenating amplicons, to generate several amplicon  
308 sequences per Nanopore read, the efficiency of this approach was modest, yielding only a two-fold increase in  
309 throughput at present. Further, the needs for concatenation and higher inputs required several PCR amplification  
310 steps that could have contributed to the observed shifts in allele frequencies leading to differing assignments on the  
311 different platforms. Turnaround time in the present study was mostly restricted by the computational capacity of the  
312 portable laptop used for the computational analysis. Specifically, base calling by translating the raw electrical signal  
313 recorded by the minION sequencer into fastq nucleotide reads proved to be the most time-consuming step, requiring  
314 up to several days in computing time.

315         However, despite the limitations associated with the Nanopore platform described above, the stock  
316 composition of coho in the Gulf of Alaska also confounded accuracy and fidelity of stock assignment. Most  
317 importantly, most salmon sampled and assessed during the Gulf of Alaska expedition were assigned to Southeastern  
318 Alaska and adjacent British Columbia coast stocks (SEAK, HecLow+HStr, SC + SFj). These stocks are poorly  
319 represented in the queried baseline and stocks from northern Alaska are very sparse so that fish from such origin  
320 often get assigned to the SEAK with poor confidence. This meant that even on the Ion Torrent platform, assignment  
321 probabilities were low, causing small differences in SNP content between the two platforms to lead to alternating  
322 assignment between these stocks (i.e. SEAK assignment on Nanaopore being assigned to HecLow+HStr and SC +  
323 SFj on Ion Torrent). Indeed, stock assignment on the Ion Torrent platform using an updated and expanded baseline  
324 and primer set, resulted in high confidence assignment of many of these individuals to Kynoch and Mussel Inlets, a  
325 spatially close reporting unit on the Northern BC coast that was poorly represented in the original baseline (C.  
326 Neville, personal communication). This suggests that new SNP loci included in the updated primer set and baseline  
327 were able to resolve these stocks at higher confidence and assign them to the appropriate stock (Beacham et al.  
328 2020). Fortunately, all of the current limitations mentioned above can be addressed in further development and we  
329 expect significant improvements in all fields, ultimately delivering a high throughput, real-time, in-field sequencing  
330 platform.

331 Advances to the Nanopore platform, sample preparation, as well as computational  
332 infrastructure will improve turnaround, throughput, and fidelity

333 While we were successful in providing a proof-of-principle study demonstrating that the Nanopore platform  
334 is capable of in-field genotyping, the throughput, fidelity, and turnaround, remained below the level needed to put  
335 this platform into standard operation for GSI by SNP genotyping. Several modifications in the workflow are planned  
336 to improve the throughput. Currently, barcoding relies on inefficient blunt-end ligation of the barcoding adapters to  
337 the PCR amplicons, leading to up to 50% unbarcoded amplicons and therefore wasting a large portion of sequencing  
338 capacity. Including the ligation adapter sequences needed to add the barcodes in the PCR primers will improve the  
339 efficacy of barcoding by circumventing the inefficient and laborious blunt-end ligation. This will improve  
340 sequencing throughput, while at the same time speeding up the sample preparation by approximately one hour. Next,  
341 concatenation efficiency is currently relatively low, increasing throughput only two-fold. While large concatemers  
342 approaching 10kb were observed, they were relatively rare. Optimized concatenation conditions by adjusting the  
343 reaction conditions such as annealing temperature and duration should exponentially improve throughput by both  
344 increasing the relative abundance of concatenated amplicons, as well as the total length of concatemers. Further  
345 workflow improvements could include pre-aliquoting of DNA extraction solution, barcodes, and primers, as well as  
346 bead cleaning materials in 96 well plates before heading into the field, which should reduce an additional two hours  
347 of sample preparation, as well as reduce the risk of cross-contamination in the field. Together, these improvements  
348 should bring the total sample preparation time to about 10h, with approximately half the time being hands-on.

349 The major current bottleneck in turnaround time is the time that base calling takes on the portable laptop  
350 computer used in the present study. GPU-enabled basecalling, like the Nanopore computation unit minIT, can  
351 provide real-time base calling to fastq and is currently being tested in the follow-up work to the present study. Actual  
352 real-time basecalling will bring the workflow in the neighbourhood of the desired 24h turnaround time.

353 An additional issue for using Nanopore sequencing is the low accuracy of the sequencing platform at the  
354 time of this project using the R9 flow cells. This low accuracy requires excessively high alignment coverage at SNP  
355 locations to ensure accurate SNP calling. However, newer Nanopore flow cells promise greatly increased accuracy  
356 (e.g., 99.999% for R10) due to “a longer barrel and dual reader head” and have recently become available. This

357 updated flow cell technology is therefore expected to greatly improve sequencing accuracy and possibly allow the  
358 lowering of alignment thresholds for SNP calling, thereby increasing the throughput more than twofold.  
359 Improvements to the SNP calling pipeline, might enable the identification and exclusion of erroneous SNP calls due  
360 to the ability to calculate the p-error associated with SNP calls, thereby increasing accuracy and repeatability.  
361 Finally, in selecting SNP loci for inclusion in GSI baselines, consideration of the types of sequences that are most  
362 problematic for Nanopore sequencing (e.g. homopolymer tracts) could go a long way to improving performance  
363 across platforms. Testing power in coastwide baselines once these problematic loci are excluded will be an important  
364 future step. Extrapolating the above-mentioned improvements would improve the current throughput of 96  
365 individuals per flow cell by more than an order of magnitude, thereby enabling cost-effective real-time and/or field-  
366 based application of the platform.

367 Currently, Nanopore-based SNP GSI is an experimental in-field stock identification tool. Turnaround of  
368 several days and throughput limited to only 96 individuals per flow cell limit its attractiveness for a wider user base.  
369 Future improvements of the sequencing platform, the sample preparation procedure, as well as the computational  
370 infrastructure will greatly improve throughput and turnaround for this. This should enable the application of  
371 Nanopore-based SNP GSI for near-real-time stock management of variable batch sizes at-sea or in remote locations.  
372 Further, parallel sequencing on several flow cells using the Oxford Nanopore GridION, which can employ five flow  
373 cells simultaneously, would enable dynamic real-time stock identification using variable batch sizes from dozens to  
374 hundreds of individuals. In the event that rapid turnaround is required, the sequencing library can also be spread  
375 across several flow cells on the GridION. Together, these updates would greatly improve the abilities of multiple  
376 user groups including government, Indigenous communities, and conservation organizations to conduct GSI for  
377 safeguarding populations at risk, while allowing sustainable harvest of healthy populations.

## 378 Acknowledgements

379 The authors would like to thank the following individuals for their contribution to the expedition and to the  
380 manuscript: Richard Beamish, Brian Riddell, and the NPAFC secretariat for the organization of the 2019 Gulf of  
381 Alaska expedition. The entire scientific crew of the 2019 GoA expedition: Evgeny Pakhomov, Gerard Foley, Brian  
382 P.V. Hunt, Arkadii Ivanov, Hae Kun Jung, Gennady Kantakov, Anton Khleborodov, Chrys Neville, Vladimir

383 Radchenko, Igor Shurpa, Alexander Slabinsky, Shigehiko Urawa, Anna Vazhova, Vishnu Suseelan , Charles Waters,  
 384 Laurie Weitkamp, and Mikhail Zuev. The crew of the research vessel Professor Kaganovskiy. Charlie Waters for  
 385 providing an R script for catch visualization. Chrys Neville for the contribution of catch data. This research was  
 386 supported by Pacific Salmon Commission, Pacific Salmon Foundation, and Fisheries and Oceans Canada and the  
 387 Canadian Coast Guard (DFO CCG). CMD was supported by a fellowship through the Pacific Salmon Foundation  
 388 and MITACS.

## 389 References

- 390 Atlas, W. I., Ban, N. C., Moore, J. W., Tuohy, A. M., Greening, S., Reid, A. J., ... & Connors, K. 2021. "Indigenous  
 391 systems of management for culturally and ecologically resilient Pacific salmon (*Oncorhynchus* spp.) fisheries."  
 392 *BioScience*, 71(2), 186-204.
- 393 Beacham, Terry D., Colin G. Wallace, Kim Jonsen, Brenda McIntosh, John R. Candy, Eric B. Rondeau, Jean-  
 394 Sébastien Moore, Louis Bernatchez, and Ruth E. Withler. 2020. "Accurate Estimation of Conservation Unit  
 395 Contribution to Coho Salmon Mixed-Stock Fisheries in British Columbia, Canada Using Direct DNA  
 396 Sequencing for Single Nucleotide Polymorphisms." *Canadian Journal of Fisheries and Aquatic Sciences*.  
 397 *Journal Canadien Des Sciences Halieutiques et Aquatiques*, no. ja.  
 398 <https://www.nrcresearchpress.com/doi/abs/10.1139/cjfas-2019-0339>.
- 399 Beacham, Terry D., Colin Wallace, Cathy MacConnachie, Kim Jonsen, Brenda McIntosh, John R. Candy, Robert H.  
 400 Devlin, and Ruth E. Withler. 2017. "Population and Individual Identification of Coho Salmon in British  
 401 Columbia through Parentage-Based Tagging and Genetic Stock Identification: An Alternative to Coded-Wire  
 402 Tags." *Canadian Journal of Fisheries and Aquatic Sciences*. *Journal Canadien Des Sciences Halieutiques et*  
 403 *Aquatiques* 74 (9): 1391–1410.
- 404 Beacham, Terry D., Colin Wallace, Cathy MacConnachie, Kim Jonsen, Brenda McIntosh, John R. Candy, and Ruth  
 405 E. Withler. 2018. "Population and Individual Identification of Chinook Salmon in British Columbia through  
 406 Parentage-Based Tagging and Genetic Stock Identification with Single Nucleotide Polymorphisms." *Canadian*  
 407 *Journal of Fisheries and Aquatic Sciences*. *Journal Canadien Des Sciences Halieutiques et Aquatiques* 75 (7):  
 408 1096–1105.
- 409 Campbell, Nathan R., Stephanie A. Harmon, and Shawn R. Narum. 2015. "Genotyping-in-Thousands by Sequencing  
 410 (GT-Seq): A Cost Effective SNP Genotyping Method Based on Custom Amplicon Sequencing." *Molecular*  
 411 *Ecology Resources* 15 (4): 855–67.
- 412 Cederholm, C. Jeff, Matt D. Kunze, Takeshi Murota, and Atuhiro Sibatani. 1999. "Pacific Salmon Carcasses:  
 413 Essential Contributions of Nutrients and Energy for Aquatic and Terrestrial Ecosystems." *Fisheries* 24 (10): 6–  
 414 15.
- 415 Cook, Rodney C., and I. Guthrie. 1987. "In-Season Stock Identification of Sockeye Salmon (*Oncorhynchus Nerka*)  
 416 Using Scale Pattern Recognition." *Canadian Special Publication of Fisheries and Aquatic sciences/Publication*  
 417 *Speciale Canadienne Des Sciences Halieutiques et Aquatiques* 96: 327–34.
- 418 Cornelis, Senne, Yannick Gansemans, Lieselot Deleye, Dieter Deforce, and Filip Van Nieuwerburgh. 2017.  
 419 "Forensic SNP Genotyping Using Nanopore MinION Sequencing." *Scientific Reports* 7 (February): 41759.
- 420 Dann, T. H., Habicht, C., Baker, T. T., & Seeb, J. E. 2013. "Exploiting genetic diversity to balance conservation and  
 421 harvest of migratory salmon". *Canadian Special Publication of Fisheries and Aquatic sciences/Publication*  
 422 *Speciale Canadienne Des Sciences Halieutiques et Aquatiques* 70(5): 785-793.
- 423 Gilbey, John, Vidar Wennevik, Ian R. Bradbury, Peder Fiske, Lars Petter Hansen, Jan Arge Jacobsen, and Ted  
 424 Potter. 2017. "Genetic Stock Identification of Atlantic Salmon Caught in the Faroese Fishery." *Fisheries*  
 425 *Research* 187 (March): 110–19.

- 426 Hinch, S. G., S. J. Cooke, A. P. Farrell, K. M. Miller, M. Lapointe, and D. A. Patterson. 2012. "Dead Fish  
427 Swimming: A Review of Research on the Early Migration and High Premature Mortality in Adult Fraser River  
428 Sockeye Salmon *Oncorhynchus Nerka*." *Journal of Fish Biology* 81 (2): 576–99.
- 429 Jefferts, K. B., P. K. Bergman, and H. F. Fiscus. 1963. "A Coded Wire Identification System for Macro-Organisms."  
430 *Nature* 198 (4879): 460–62.
- 431 Lichatowich, Jim. 2001. *Salmon Without Rivers: A History Of The Pacific Salmon Crisis*. Island Press.
- 432 Li, Heng, and Richard Durbin. 2009. "Fast and Accurate Short Read Alignment with Burrows–Wheeler Transform."  
433 *Bioinformatics* 25 (14): 1754–60.
- 434 Li, Heng, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis,  
435 Richard Durbin, and 1000 Genome Project Data Processing Subgroup. 2009. "The Sequence Alignment/Map  
436 Format and SAMtools." *Bioinformatics* 25 (16): 2078–79.
- 437 Mikheyev, Alexander S., and Mandy M. Y. Tin. 2014. "A First Look at the Oxford Nanopore MinION Sequencer."  
438 *Molecular Ecology Resources* 14 (6): 1097–1102.
- 439 Miller, Kristina M., Amy Teffer, Strahan Tucker, Shaorong Li, Angela D. Schulze, Marc Trudel, Francis Juanes, et  
440 al. 2014. "Infectious Disease, Shifting Climates, and Opportunistic Predators: Cumulative Factors Potentially  
441 Impacting Wild Salmon Declines." *Evolutionary Applications* 7 (7): 812–55.
- 442 Miller, Kristina M., Ruth E. Withler, and Terry D. Beacham. 1996. "Stock Identification of Coho Salmon  
443 (*Oncorhynchus kisutch*) Using Minisatellite DNA Variation." *Canadian Journal of Fisheries and Aquatic  
444 Sciences. Journal Canadien Des Sciences Halieutiques et Aquatiques* 53 (1): 181–95.
- 445 Moran, Benjamin M., and Eric C. Anderson. 2019. "Bayesian Inference from the Conditional Genetic Stock  
446 Identification Model." *Canadian Journal of Fisheries and Aquatic Sciences. Journal Canadien Des Sciences  
447 Halieutiques et Aquatiques* 76 (4): 551–60.
- 448 Ozerov, Mikhail, Anti Vasemägi, Vidar Wennevik, Rogelio Diaz-Fernandez, Matthew Kent, John Gilbey, Sergey  
449 Prusov, Eero Niemelä, and Juha-Pekka Vähä. 2013. "Finding Markers That Make a Difference: DNA Pooling  
450 and SNP-Arrays Identify Population Informative Markers for Genetic Stock Identification." *PloS One* 8 (12):  
451 e82434.
- 452 Paradis, Emmanuel, and Klaus Schliep. 2019. "Ape 5.0: An Environment for Modern Phylogenetics and  
453 Evolutionary Analyses in R." *Bioinformatics* 35 (3): 526–28.
- 454 Quick, Joshua, Nicholas J. Loman, Sophie Duraffour, Jared T. Simpson, Ettore Severi, Lauren Cowley, Joseph Akoi  
455 Bore, et al. 2016. "Real-Time, Portable Genome Sequencing for Ebola Surveillance." *Nature* 530 (7589): 228–  
456 32.
- 457 Robinson, James T., Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S. Lander, Gad Getz, and Jill  
458 P. Mesirov. 2011. "Integrative Genomics Viewer." *Nature Biotechnology* 29 (1): 24–26.
- 459 Schlecht, Ulrich, Janine Mok, Carolina Dallett, and Jan Berka. 2017. "ConcatSeq: A Method for Increasing  
460 Throughput of Single Molecule Sequencing by Concatenating Short DNA Fragments." *Scientific Reports* 7 (1):  
461 5252.
- 462 Winans, G. A., Aebersold, P. B., Urawa, S., and Varnavskaya, N. V. 1994. "Determining continent of origin of chum  
463 salmon (*Oncorhynchus keta*) using genetic stock identification techniques: status of allozyme baseline in Asia."  
464 *Canadian Journal of Fisheries and Aquatic Sciences. Journal Canadien Des Sciences Halieutiques et  
465 Aquatiques*, 51(S1), 95-113.
- 466 Wood, Chris C., Dennis T. Rutherford, and Skip McKinnell. 1989. "Identification of Sockeye Salmon  
467 (*Oncorhynchus Nerka*) Stocks in Mixed-Stock Fisheries in British Columbia and Southeast Alaska Using  
468 Biological Markers." *Canadian Journal of Fisheries and Aquatic Sciences. Journal Canadien Des Sciences  
469 Halieutiques et Aquatiques* 46 (12): 2108–20.
- 470

## 471 Data accessibility and benefit-sharing statement

### 472 Data accessibility statement

473 **Data analysis pipeline:** The full pipeline to genotype salmon from nanopore data titled “nano2geno” (n2g) can be  
474 found at <https://github.com/bensutherland/nano2geno/>.

475 **Primer and genotype information:** Primer sequences and genotype information have previously been published by  
476 Beacham et al. (Beacham et al. 2017; Appendix A2).

477 **Genetic Data:** All raw nanopore sequence reads analyzed in this paper are deposited in the SRA under BioProject:  
478 PRJNA796718 (SRR17593964 - SRR17593966).

479 **Sample metadata:** Metadata on the individuals in this study is also stored associated with BioProject:  
480 PRJNA796718 under the BioSamples SAMN24907542-SAMN24907622.

481 **Genotype baseline data:** The genotype baseline used for stock identification with rubias in this manuscript is based  
482 on Beacham et al. 2017 and 2020 and is available on DataDryad (<https://doi.org/10.5061/dryad.g4f4qrfs3>)

### 483 Benefit sharing statement

484 Benefits Generated: Benefits from this research accrue from the sharing of our methodology and reference data as  
485 described throughout the manuscript and available under the repositories mentioned in data accessibility statement.

## 486 Author Contributions

487 C.M. Deeg, B.J. G. Sutherland, and K.M. Miller designed research. C.M. Deeg performed research. T.J. Ming, C.  
488 Wallace, K. Jonsen, K.L. Flynn, E.B. Rondeau, and T.D. Beacham contributed new reagents or analytical tools. C.M.  
489 Deeg, B.J. G. Sutherland, and E.B. Rondeau, analyzed data. C.M. Deeg, B.J. G. Sutherland, and K.M. Miller wrote  
490 the paper.

491

## 492 Tables and Figures

493 Table1: Relative proportion of top reporting units (contribution >3%) to the overall mixture of coho salmon. Only  
494 individuals that had successful stock ID on all three GSI runs are included. Reporting Units: SEAK: Southeast  
495 Alaska; LSTK: Lower Stikine River; NCS: North Coast Streams (BC); HecLow+HStr: Lower Hecate Strait and  
496 Haro Strait; SC + SFj: Southern Coastal Streams, Queen Charlotte Strait, Johnston Strait and Southern Fjords; CR:  
497 Columbia River; COWA: Coastal Washington; LNASS: Lower Nass River; WVI: West Vancouver Island; OR:  
498 Oregon.

499  
500 Figure 1: Workspace aboard the Professor Kaganovsky vessel during the International Year of the salmon signature  
501 expedition.

502  
503 Figure 2: Simplified wet-lab workflow for DNA extraction, amplification, barcoding, and concatenation before  
504 sequencing and pipeline of the following computational analysis. DNA is shown in black, amplification primers in  
505 green, fish ID barcodes in olive, concatenation adapters in red/blue, and sequencing adapters in purple.

506  
507 Figure 3: Number of reads per amplicon per individual (barcode) of Nanopore sequencing runs. The violin plot  
508 shows the distribution of number of reads assigned to unique SNP-containing amplicons within an individual. Green  
509 and blue colors denote the two separate sequencing runs during the IYS expedition (top), and black indicates the run  
510 at the laboratory (PBS; bottom). Above each individual violin plot is the total number of amplicons for that  
511 individual for which sufficient reads were present to call the genotype, color indicates if enough amplicons were  
512 called for downstream analysis (black) or not (red). The order of individuals is matched in the top and bottom plots.

513  
514 Figure 4: Comparison of sequence alignment of Nanopore and Ion Torrent sequences from the same individual  
515 against a SNP locus preceded by a homopolymer tract. Nanopore sequences show a higher number of indels,  
516 specifically associated with the poly-T homopolymer tract (145-151bp) directly preceding the SNP location (152bp).  
517 Alignment was visualized here using IGV (Robinson et al. 2011)

518  
519 Figure 5: Principal coordinate analysis (PCoA) of SNP calls of individuals passing threshold in all datasets. SNP  
520 calls based on Nanopore sequences generated during the IYS expedition shown in blue (“nano\_IYS”), and the same  
521 individuals reanalyzed upon return using the same workflow shown in purple (“nano\_PBS”). Ion Torrent reads  
522 scored with the n2g pipeline are shown in red (“ion\_n2g”) and scores derived from the Ion Torrent variant caller are  
523 shown in green (“ion\_vc”).

524

525 Figure 6: Relative proportion of reporting units to the overall mixture of coho salmon. Only individuals that had  
526 passed the stock ID threshold (>50% of SNPs called) on all three GSI runs are included. Reporting Units: SEAK:  
527 Southeast Alaska; LSTK: Lower Stikine River; HecLow+HStr: Lower Hecate Strait and Haro Strait; SC + SFj:  
528 Southern Coastal Streams, Queen Charlotte Strait, Johnston Strait and Southern Fjords; CR: Columbia River;  
529 COWA: Coastal Washington.

530

For Review Only

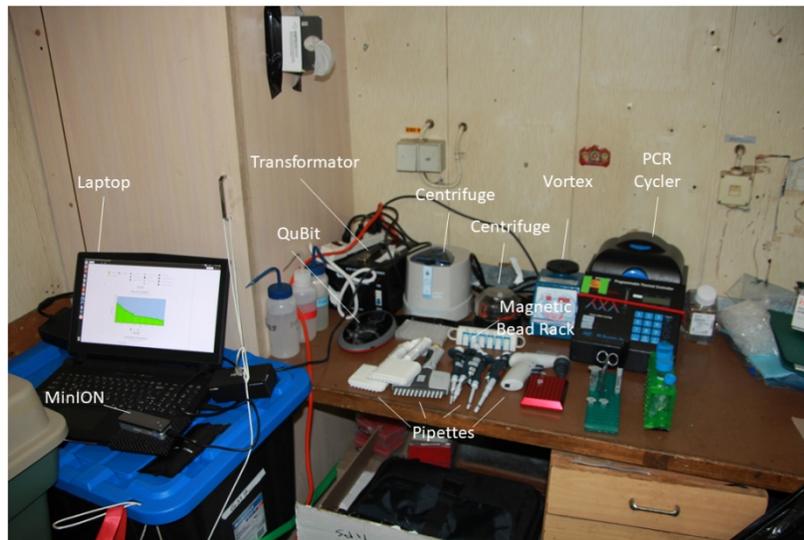


Figure 1: Workspace aboard the Professor Kaganovsky vessel during the International Year of the salmon signature expedition.

855x481mm (38 x 38 DPI)

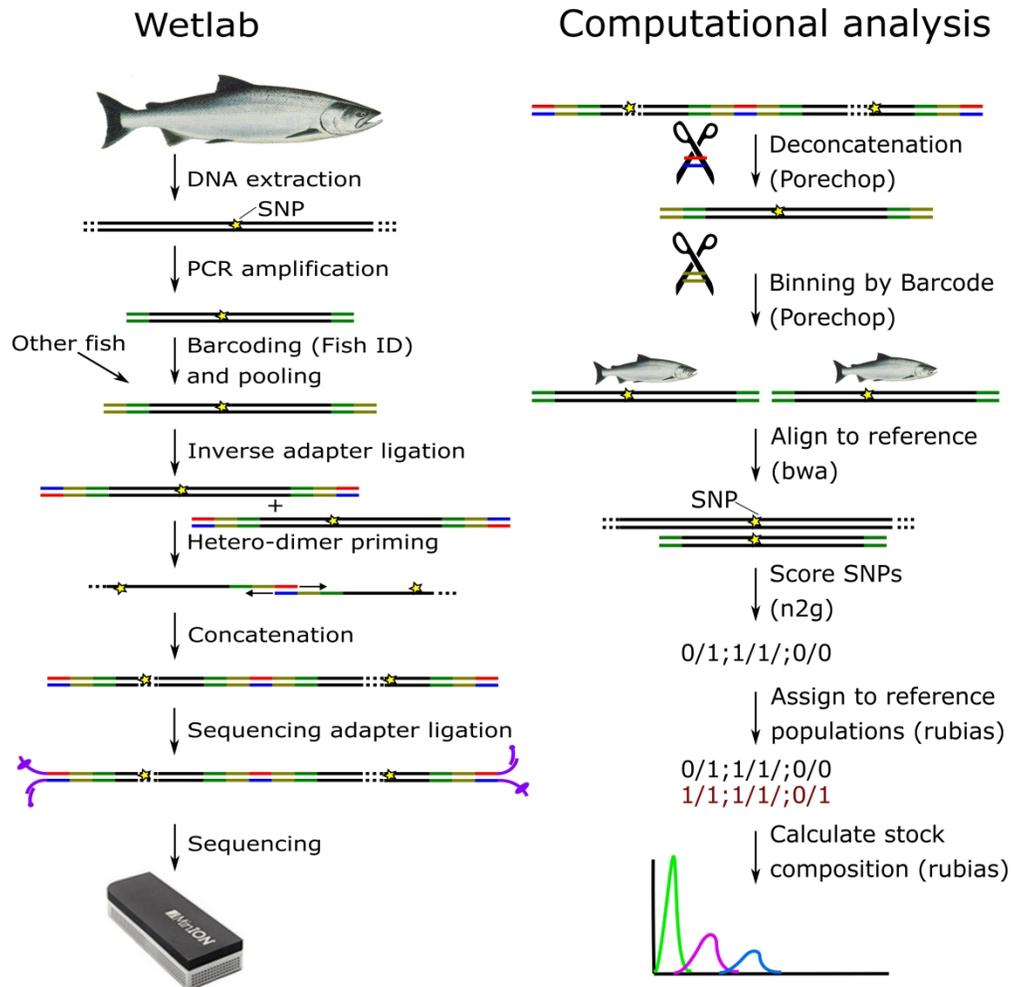


Figure 2: Simplified wet-lab workflow for DNA extraction, amplification, barcoding, and concatenation before sequencing and pipeline of the following computational analysis. DNA is shown in black, amplification primers in green, fish ID barcodes in olive, concatenation adapters in red/blue, and sequencing adapters in purple.

1263x1267mm (72 x 72 DPI)

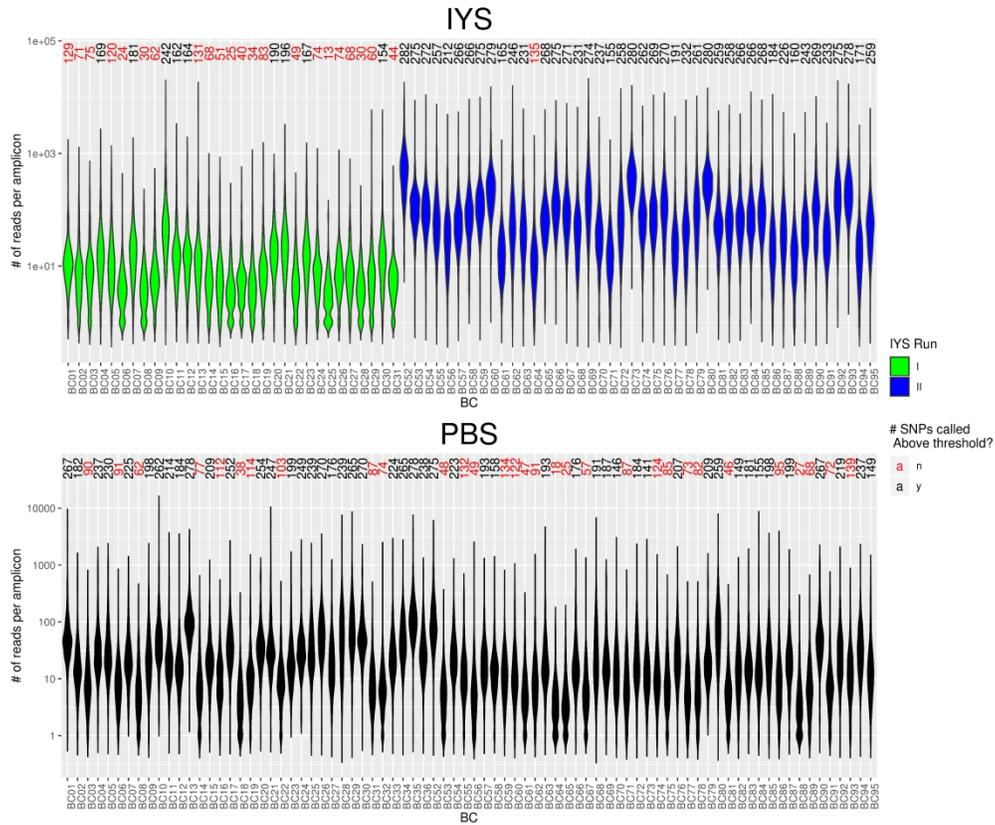


Figure 3: Number of reads per amplicon per individual (barcode) of Nanopore sequencing runs. The violin plot shows the distribution of number of reads assigned to unique SNP-containing amplicons within an individual. Green and blue colors denote the two separate sequencing runs during the IYS expedition (top), and black indicates the run at the laboratory (PBS; bottom). Above each individual violin plot is the total number of amplicons for that individual for which sufficient reads were present to call the genotype, color indicates if enough amplicons were called for downstream analysis (black) or not (red). The order of individuals is matched in the top and bottom plots.

299x250mm (300 x 300 DPI)



Figure 4: Comparison of sequence alignment of Nanopore and Ion Torrent sequences from the same individual against a SNP locus preceded by a homopolymer tract. Nanopore sequences show a higher number of indels, specifically associated with the poly-T homopolymer tract (145-151bp) directly preceding the SNP location (152bp). Alignment was visualized here using IGV (Robinson et al. 2011)

394x172mm (118 x 118 DPI)

Rank	Ion Torrent (ion_vc)			Nanopore (nano_PBS)		
	Repunit	Proportion	SD	Repunit	Proportion	SD
1	SEAK	0.437678	0.109758	SEAK	0.662083	0.218561
2	HecLow+H Str	0.178637	0.057264	LSTK	0.205116	NA
3	LSTK	0.068878	NA	CR	0.050276	0.012993
4	SC+SFj	0.067989	0.025318	COWA	0.042244	0.011583
5	CR	0.067939	0.01403			
6	NCS	0.036009	0.004052			
7	OR	0.034352	0.010704			
8	WVI	0.033487	0.009144			
9	LNASS	0.032288	0.022742			

For Review Only

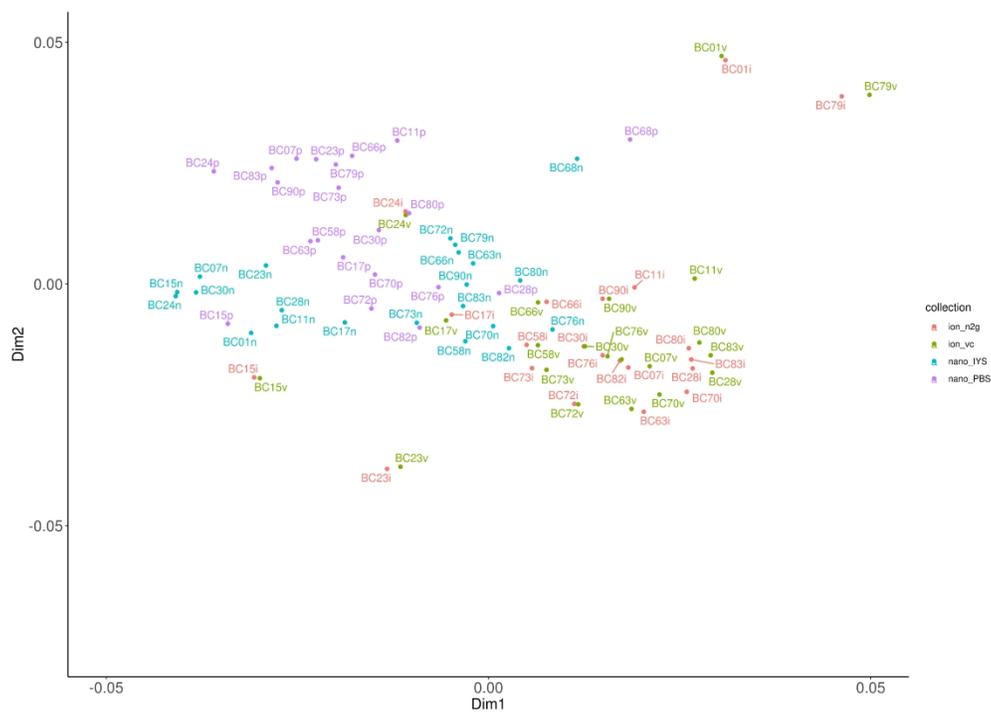


Figure 5: Principal coordinate analysis (PCoA) of SNP calls of individuals passing threshold in all datasets. SNP calls based on Nanopore sequences generated during the IYS expedition shown in blue ("nano\_IYS"), and the same individuals reanalyzed upon return using the same workflow shown in purple ("nano\_PBS"). Ion Torrent reads scored with the n2g pipeline are shown in red ("ion\_n2g") and scores derived from the Ion Torrent variant caller are shown in green ("ion\_vc").

350x250mm (257 x 257 DPI)

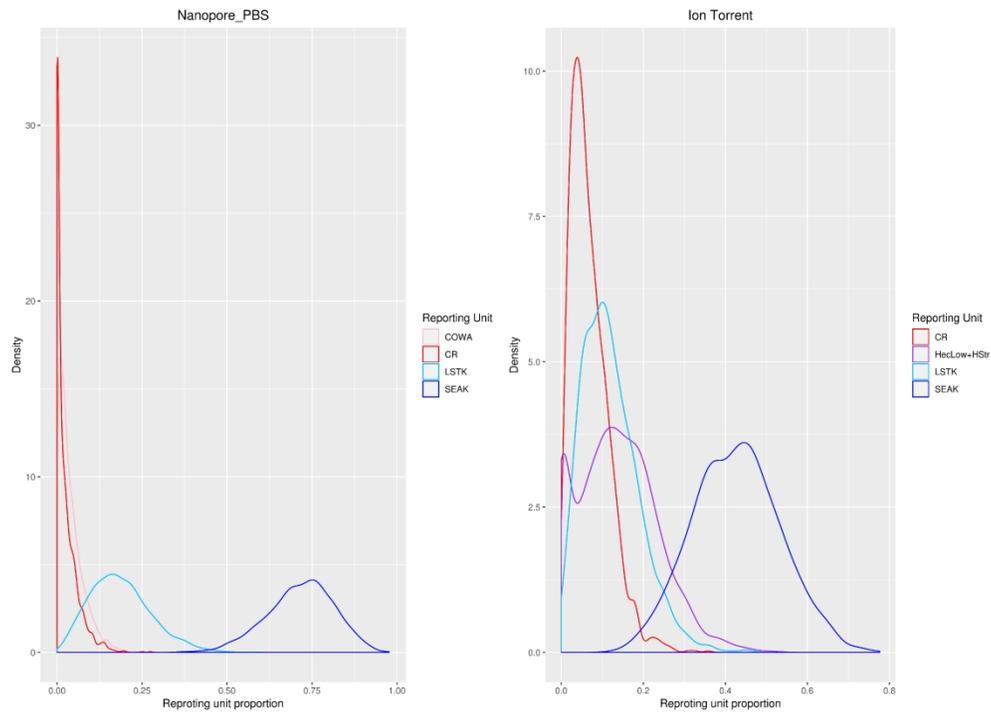


Figure 6: Relative proportion of reporting units to the overall mixture of coho salmon. Only individuals that had passed the stock ID threshold (>50% of SNPs called) on all three GSI runs are included. Reporting Units: SEAK: Southeast Alaska; LSTK: Lower Stikine River; HecLow+HStr: Lower Hecate Strait and Haro Strait; SC + SFj: Southern Coastal Streams, Queen Charlotte Strait, Johnston Strait and Southern Fjords; CR: Columbia River; COWA: Coastal Washington.

350x250mm (257 x 257 DPI)