


## RESOURCE ARTICLE

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

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. Rondeau<sup>3</sup> | Terry D. Beacham<sup>3</sup>  | Kristina M. Miller<sup>1,3</sup>

<sup>1</sup>Forest and Conservation Sciences, University of British Columbia, Vancouver, British Columbia, Canada

<sup>2</sup>Pacific Salmon Foundation, Vancouver, British Columbia, Canada

<sup>3</sup>Fisheries and Oceans Canada, Pacific Biological Station, Nanaimo, British Columbia, Canada

## Correspondence

Christoph M. Deeg, Forest and Conservation Sciences, University of British Columbia, Vancouver, BC, Canada.  
Email: [chdeeg@mail.ubc.ca](mailto:chdeeg@mail.ubc.ca)

## Funding information

Mitacs, Grant/Award Number: IT13895; Pacific Salmon Foundation; Pacific Salmon Commission; Fisheries and Oceans Canada

Handling Editor: Shotaro Hirase

## Abstract

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

## KEYWORDS

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

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial](https://creativecommons.org/licenses/by-nc/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2022 The Authors. *Molecular Ecology Resources* published by John Wiley & Sons Ltd.

## 1 | INTRODUCTION

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

To inform mixed-stock management, stock identification has in the distant past utilized characteristic scale and parasite patterns as well as the marking of hatchery-enhanced fish by coded-wire tagging (Cook & Guthrie, 1987; Jefferts et al., 1963; Wood et al., 1989). More recently, genetic stock identification (GSI) using allozyme, minisatellite, microsatellite and ultimately single nucleotide polymorphisms (SNPs) as markers has proven superior in delivering high-throughput insights into the stock composition of salmon (Beacham et al., 2017, 2018; Miller et al., 1996; Winans et al., 1994). Specifically, the large baseline of population-specific SNP frequencies and targeted amplification of such SNP loci now allow for unprecedented resolution of stock origin in many species of salmon at reduced biases (Beacham et al., 2017, 2018; Gilbey et al., 2017; Ozerov et al., 2013). However, current sequencing approaches, based on second-generation sequencing platforms (e.g., Illumina and Ion Torrent), mean that only sequencing large batches of individuals, known as “genotyping by the thousands” (GT-seq), is economically sensible (Beacham et al., 2017, 2018; Campbell et al., 2015). These approaches require a specialized laboratory and several days’ turnover for the library preparation and sequencing, even under highly automated settings. These constraints limit the utility of SNP-based GSI for real-world scenarios that are often spatially or temporally restricted, because samples need to be transported to the laboratory for analysis, as has been the case for most GSI methods to date. Specifically, for time-sensitive stock-specific harvest management decisions, an in-field real-time SNP-based GSI approach with greater flexibility in sample batch size would be desirable.

Recent advances in third-generation single-molecule sequencing platforms such as the Oxford Nanopore minION allow real-time sequencing on a pocket-sized portable sequencer that requires little library preparation, therefore enabling sequencing in remote

locations (Mikheyev & Tin, 2014; Quick et al., 2016). However, several technical hurdles to adapting Nanopore sequencing to SNP GSI exist. While Nanopore sequencing can yield extremely long reads, the number of sequencing pores and their loading rate is limited, resulting in low throughput when sequencing short reads such as amplicons. An additional problem is the relatively high error rate inherent to this novel technology. Since the SNP GSI protocols are based on the amplification of short amplicons via targeted multiplex PCR, sequencing throughput of such short amplicons on the Nanopore platform is comparatively low, as the number of sequencing pores is the rate-limiting factor. This is especially problematic because high coverage is needed to compensate for the higher error rate of Nanopore-generated sequences. A promising approach to overcome these limitations is the concatenation of PCR amplicons that allows the sequencing of several amplicons within a single read, thereby exponentially increasing throughput for genotyping (Cornelis et al., 2017; Schlecht et al., 2017).

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

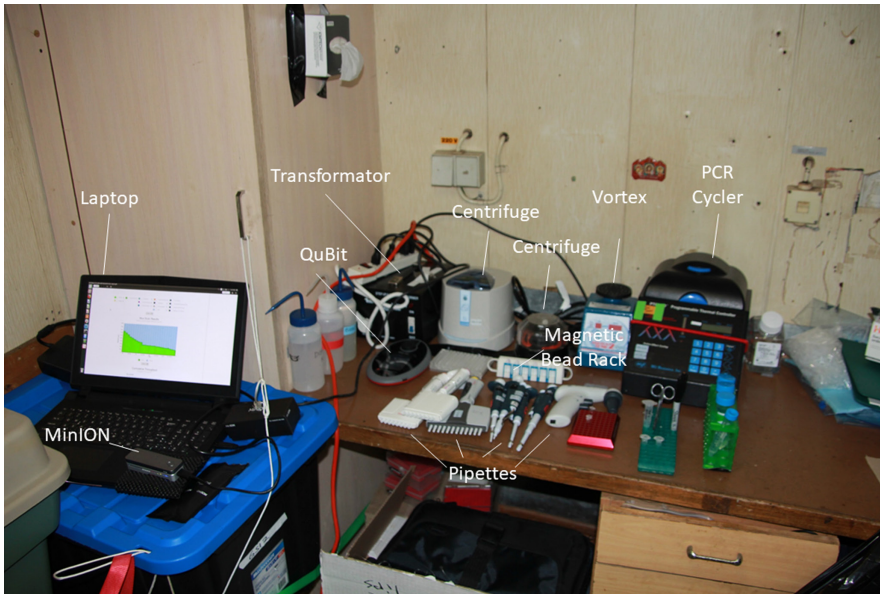
## 2 | MATERIALS AND METHODS

### 2.1 | Field laboratory equipment and workspace

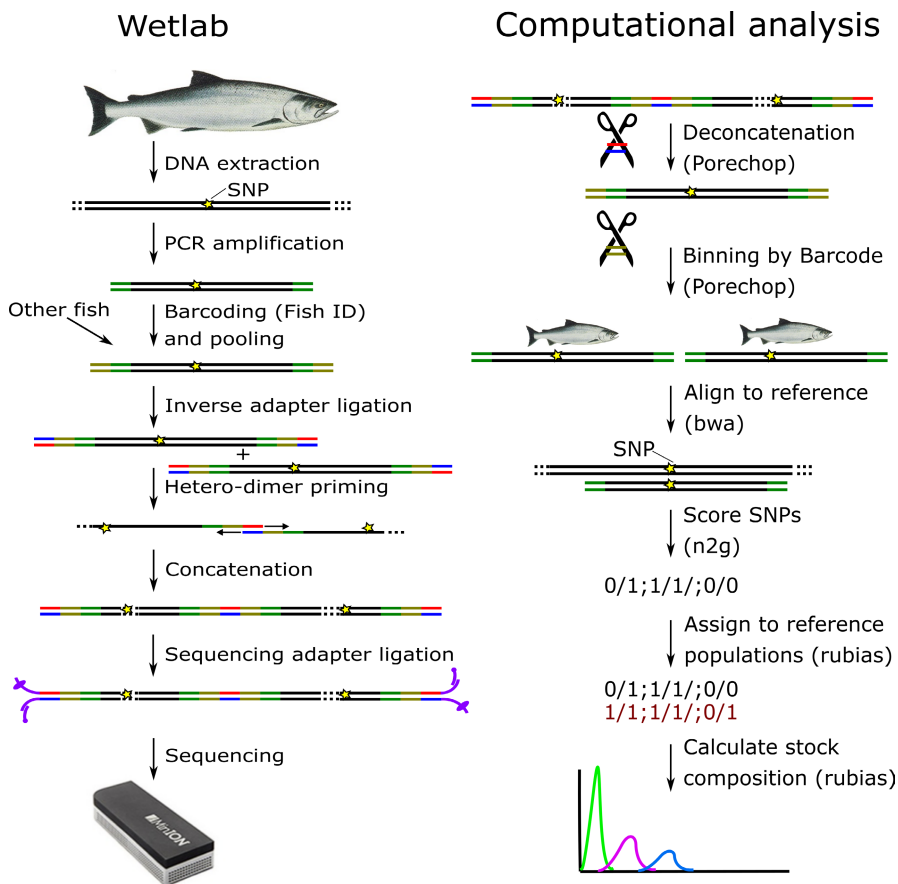
The field equipment onboard the *Professor Kaganovsky* research trawler consisted of a PCR thermocycler, a miniplate centrifuge, a microcentrifuge, a Qubit fluorimeter (Thermo Fisher), a vortexer, a minION sequencer, a laptop with an Ubuntu operating system (Ubuntu version 14.06), as well as assorted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

### 2.2 | Tissue sample collection and DNA extraction

Salmon were captured by the research trawler *Professor Kaganovsky* during the 2019 IYS Signature expedition in the Gulf of Alaska (Figure S1). We collected fin clips of coho salmon (*Oncorhynchus kisutch*) and froze them individually until DNA extraction, or immediately processed once a suitable batch size had been accumulated. DNA extraction from 2 × 2 × 2-mm fin-tissue clips was performed in a 96-well PCR plate using 100 µl of QuickExtract solution (Lucigen) according to the manufacturer's instructions.



**FIGURE 1** Workspace aboard the *Professor Kaganovsky* vessel during the International Year of the salmon signature expedition



**FIGURE 2** Simplified wet-laboratory 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

### 2.3 | Multiplex PCR and barcoding

Multiplex PCR with a custom panel of primers targeting 299 loci of known SNPs was performed using 0.25  $\mu$ l of DNA extract as template using the AgriSeq HTS Library Kit Amplification Mix PCR mastermix (ThermoFisher) in a 10- $\mu$ l reaction according to Beacham et al. (2017; see Appendix A2). Primer sets targeting

multinucleotide polymorphisms (MNPs) were included in the primer panel by Beacham et al. (2017) but were excluded from the analysis (Table S1). Next, we prepared amplicons for ligation by end-prepping amplified strands with AgriSeq HTS Library Kit Pre-ligation Enzyme. ONT barcode adapters (PCR Barcoding Expansion 1–96, EXP-PBC096; Oxford Nanopore Technologies) were then ligated to the amplicons by blunt-end ligation with the Barcoding

Enzyme/Buffer of the AgriSeq HTS Library Kit according to the manufacturer's instructions. After bead-cleanup (1.2:1 bead/sample, AMPure XP beads; Beckman Coulter) we added the ligation products, barcodes and barcoding adapters (PCR Barcoding Expansion 1-96, EXP-PBC096; Oxford Nanopore Technologies) by PCR using Q5 polymerase mastermix (NEB) for individual fish identification according to the manufacturer's protocol in a 25- $\mu$ l reaction (98°C for 3 min; 25 cycles of 98°C for 10 s, 70°C for 10 s, 72°C for 25 s; 72°C for 2 min). Barcoded libraries were then pooled and cleaned using 1.2:1 bead cleanup, before DNA yield of a subset of samples (12.5%) was analysed by Qubit (dsDNA HS Assay Kit; ThermoFisher).

## 2.4 | Amplicon concatenation

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

## 2.5 | Library preparation and sequencing

The concatenated amplicons were prepared for Nanopore sequencing using the ONT Ligation Sequencing Kit (LSK109) according to

the manufacturer's instructions. In brief, after end-preparation using the Ultra II Endprep Module and bead cleanup, we ligated proprietary ONT sequencing adapters onto the concatenation adapters by blunt-end ligation using the proprietary ONT Buffer and the TA quick ligase (NEB; note: this standard sequencing step is not shown in Figure 2). After additional bead-cleanup and washing with the short fragment buffer (SFB; ONT) according to the manufacturer's protocol, we loaded the library onto a freshly primed flow cell (MIN 106 R9.4.1; ONT) according to the manufacturer's instructions.

## 2.6 | Nanopore sequencing, deconcatenation and binning

After flow cell priming and loading of the library, the flow cell was placed on the minION sequencer. Sequencing and basecalling into fast5 and fastq was performed simultaneously using MINKNOW (version 3.1.8) on an Ubuntu 14.06 platform. First, all fastq raw reads that passed default quality control in MINKNOW were combined into bins of 500,000 reads each. This had empirically been determined to be the maximum number of reads allowing simultaneous processing in the downstream analysis on our platform (Ubuntu 14.06, 31.2 GiB RAM 7700K CPU @ 4.20 GHz x8). Reads containing concatenated amplicons were deconcatenated and the concatenation adapter sequence was trimmed off the remaining sequence using PORECHOP (<https://github.com/rrwick/Porechop>) with a custom adapter file ("adapters.py") that only contained the concatenation adapter under the following settings: porechop-runner.py -i input\_raw\_reads.fastq -o output/dir -t 16 --middle\_threshold 75 --min\_split\_read\_size 100 --extra\_middle\_trim\_bad\_side 0 --extra\_middle\_trim\_good\_side 0.

We binned the deconcatenated reads by barcode corresponding to fish individuals by using PORECHOP with the provided default adapters file and the following settings: porechop-runner.py -i input\_deconcatenated\_reads.fastq -b binning/dir -t 16 --adapter\_threshold 90 --end\_threshold 75 --check\_reads 100000.

After this step, all reads from the corresponding barcode bins corresponding to the same individual across the different 500,000 sub-bins were combined for downstream analysis. See <https://github.com/bensutherland/nano2geno/> for source scripts for analysis.

## 2.7 | Alignment and SNP calling

We aligned the binned reads to the reference amplicon sequences described by Beacham et al. (2017) using BWA-MEM and indexed them using SAMTOOLS (Beacham et al., 2017; Li & Durbin, 2009; Li et al., 2009). Alignment statistics for all loci were generated using PYSAMSTATS (<https://github.com/alimanfoo/pysamstats>; flags: -t variation -f) and we extracted the nucleotides observed at the relevant SNP hotspot loci from the resulting file using a custom R script by looping through the results file guided by an SNP location file. Finally, we compared

the observed nucleotide distributions at SNP hotspots to the hotspot reference and variant nucleotides and scored as homozygous reference when  $\geq 66\%$  of the nucleotides were the reference allele, heterozygous when the reference allele was present  $< 66\%$  and the variant allele  $> 33\%$ , or as homozygous variant when the nucleotides were  $\geq 66\%$  the variant allele, using a custom R script to generate a numerical locus table. We visually inspected alignments determined to be problematic using the IGV VIEWER (Robinson et al., 2011). The full pipeline entitled “nano2geno” (n2g) including all custom scripts can be found at <https://github.com/bensutherland/nano2geno/> (Figure 2).

## 2.8 | Mixed-stock analysis

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

## 2.9 | Ion torrent sequencing

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

## 2.10 | Concordance assessment

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

## 3 | RESULTS

### 3.1 | In-field Nanopore sequencing

During the IYS Signature expedition to the Gulf of Alaska in February and March 2019, in-field SNP GSI was performed on coho salmon as the tissues became available. A total of 75 coho salmon were analysed in two sequencing runs at different points during the expedition, representing 77% of all coho salmon captured during the expedition.

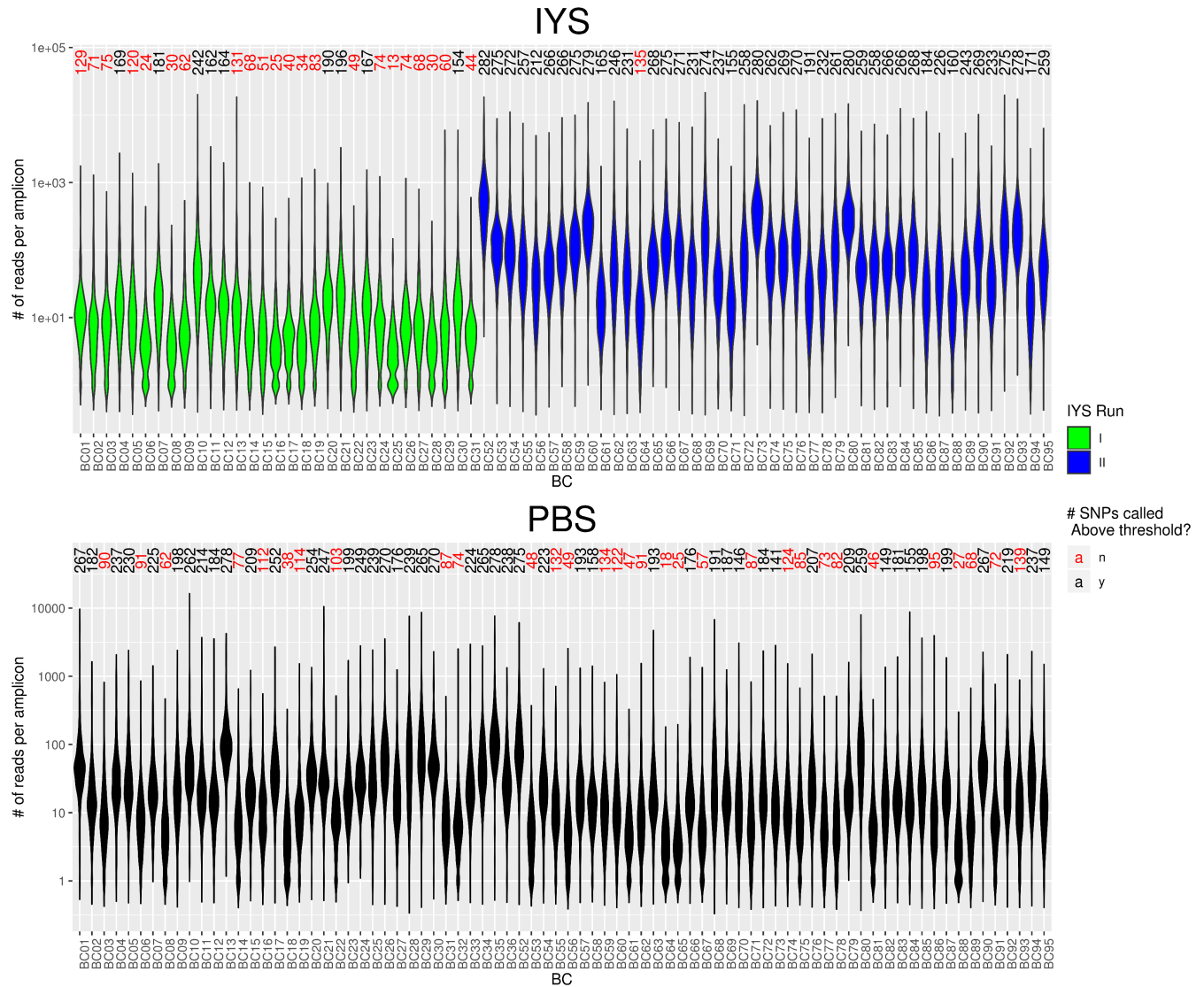
The first sequencing run was performed on February 26 and included 31 individuals. Library preparation onboard the vessel took 14 h. However, faulty flow cell priming resulted in only approximately half the detected pores being active (843 pores). Of these pores, no more than 25% were actively sequencing at any time, highlighting the challenges of utilizing sensitive equipment under field conditions including excessive ship movement. Accordingly, sequencing for 30 h and base-calling for 34 h resulted in only 1.44 million (M) reads, 49% of which passed quality control. The read length distribution showed several large, concatenated amplicons up to 7,095 bp with a mean length of 825 bp (Figure S2). Deconcatenation resulted in a read inflation by a factor of 2 $\times$  (702 thousand [k] to 1,444k reads). After binning, reads per individual ranged from 1983 to 86,467 with a mean of 13,709 reads (SD: 15,370), and 722,174 reads that were not able to be assigned (50% of the total deconcatenated reads) (Figure 3; Figures S2 and S3).

The second sequencing run was performed on March 10, 2019, with 44 coho salmon. Library preparation again took 14 h and sequencing on a new flow cell took 15 h, starting with 1502 available pores, and up to 65% actively sequencing pores, and resulted in 4.48 M reads, 76% of which passed quality control. Read lengths averaged 810 bp with a maximum length of 8023 bp (Figure S2). Due to the large number of reads and the limited power of the computer being used for the analysis, base-calling into fastq took 3 days. Deconcatenation resulted in a read inflation of a factor of 1.7 $\times$  (3.4 M to 5.8 M) (Figure S2). Reads per individual showed a mean of 67,636 (SD: 59,393; min: 11,684; max: 335,348), with 722,179 reads remaining unassigned (12%) (Figure 3; Figures S2 and S3).

Upon return from the expedition, we sequenced 80 individuals, including all those previously genotyped aboard the vessel, in a single MinION run using the expedition setup starting from the frozen tissues from the expedition. We sequenced for 42 h to maximize the total number of reads with 60% of 2048 available pores actively sequencing resulting in 5.32 M reads. Of these reads, 3.20 M passed quality control. Again, large, concatenated amplicons up to 9449 kb were observed, with a mean read length of 840 bp, and deconcatenation resulted in 4.54 M reads (1.4 $\times$  inflation) (Figure S2). The mean number of reads per bin was 29,439 (SD: 25,000) and ranged from 2969 to 128,718 reads per individual, with 1,413,626 unassigned reads (31%) (Figure 3; Figure S2).

Despite the absence of normalization between samples prior to multiplex PCR, barcoding and loading, the binning distribution





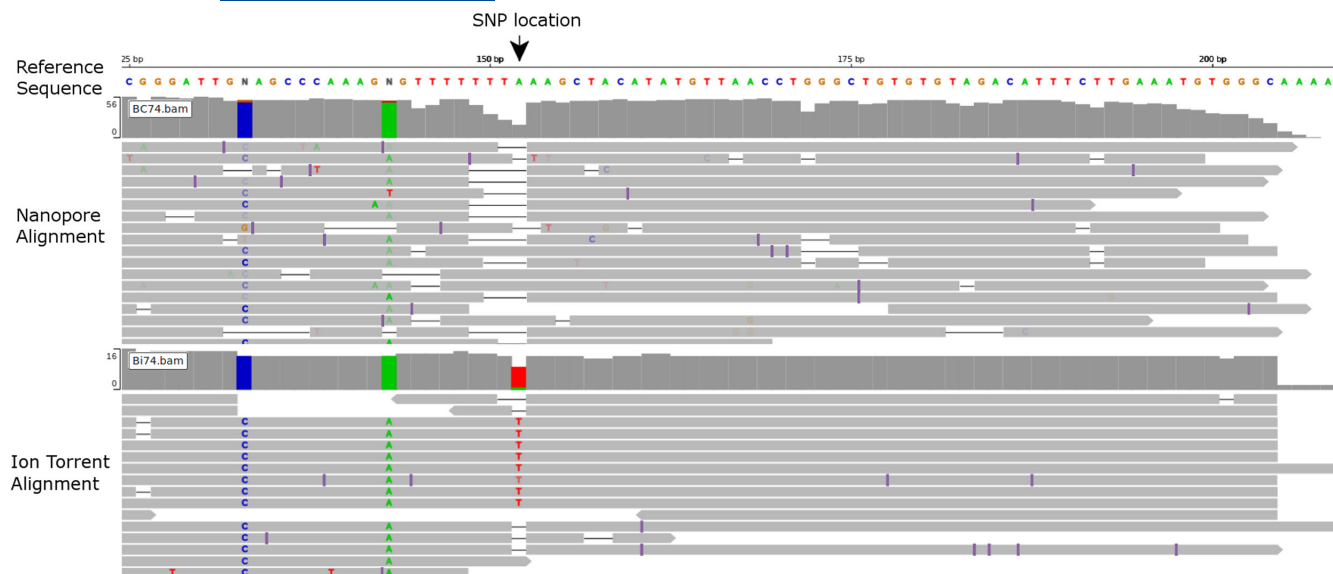
**FIGURE 3** Number of reads per amplicon per individual (barcode) of Nanopore sequencing runs. The violin plot shows the distribution of number of reads assigned to unique SNP-containing amplicons within an individual. Green and blue colours denote the two separate sequencing runs during the IYS expedition (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, and colour 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

across samples was relatively even with only a few apparent outliers observed (Figure 3; Figure S3). The minimum number of reads per individual sample necessary to cover sufficient loci (at a minimum depth of 10 sequences per locus) for downstream stock assignments (i.e., at least 141 loci per sample) is around 2000 reads (Figure 3; Figure S3).

### 3.2 | Nanopore sequencing data require loci reassessment for efficient SNP calling

After alignment to the reference sequences for SNP calling, Nanopore sequence data showed a comparatively higher error rate than Ion Torrent reads, as expected, with abundant indels that frequently led to lower alignment scores than those obtained by the

Ion Torrent data (Ion Torrent average alignment score: 25.6 MAPQ; Nanopore average alignment score: 13.9 MAPQ). Specifically, regions containing homopolymer tracts were poorly resolved, as had previously been reported (Cornelis et al., 2017). Several instances could be identified where the homopolymer presence near the SNP locus caused problematic alignments and therefore resulted in SNP calls not matching those found by the Ion Torrent on the same individual (Figure 4). Accordingly, six such loci were excluded from downstream analysis (Table S1). Other loci were excluded from the analysis due to absence of coverage (four loci) or the inability of the custom n2g pipeline to call MNPs or deletions (seven loci), bringing the number of accessed loci from 299 to 282. Other loci showing apparent differences between Nanopore and Ion Torrent sequence data ( $n = 21$ ) were retained as no apparent explanation for the discrepancies could be identified.



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

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

### 3.3 | Platform biases lead to moderately altered SNP calling compared to Ion Torrent sequencing

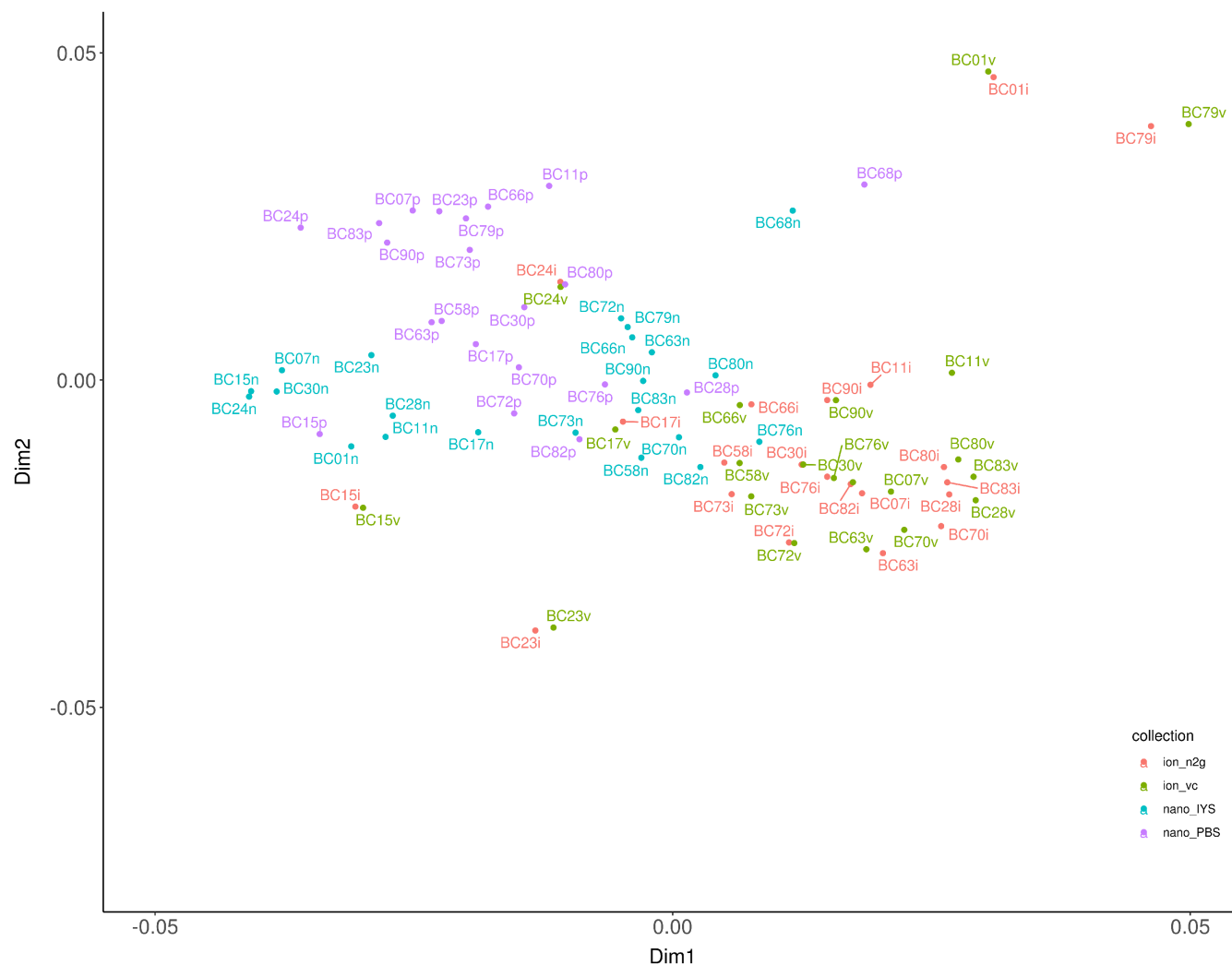
To assess the discrepancies between sequencing platforms, individuals that passed the genotyping rate threshold of 141 called loci (50% genotyping rate) in all data sets (i.e., Nanopore data during the expedition analysed with n2g: “nano IYS,” Nanopore acquired during the repeat run upon return from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

agreement on SNP call between both Nanopore runs (comparing reference or alternate scores for both alleles from nano IYS vs. nano PBS) was 84.4%, highlighting the inter-run variability associated with current Nanopore sequencing. There was a slight correlation observed between the number of Nanopore reads per individual and the concordance with Ion Torrent SNP calls, suggesting that read depth is only a minor factor influencing SNP call concordance at the current threshold of a minimal alignment depth of 10× per site for Nanopore reads (Figure S4). Excluding MNPs, deletions and homopolymer issues, the influence of the SNP calling pipeline (n2g vs. variant caller) appears negligible compared to the differences by sequencing platform (Figure 5). Accordingly, SNPs scored based on the same Ion Torrent data sequence matched in 99.21% of cases between the two genotyping pipelines.

### 3.4 | Stock assignment based on Nanopore data is moderately repeatable and differs inherently from Ion Torrent-based assignments in a subset of individuals

Stock assignment by rubias showed discrepancies between the Nanopore- and Ion Torrent-based data sets. In only 61.5% of cases did Nanopore sequences (PBS run) lead to the same top reporting unit (repunit; large-scale geographical areas such as Westcoast Vancouver Island or Lower Fraser River) assignment for individual stock ID as the Ion Torrent-based sequences (Figure 6, Table 1). Specifically, Nanopore-based repunit assignment showed higher proportions of assignments to Southeastern Alaska (SEAK) than Ion Torrent-based assignments (Figure 6, Table 1).

Nevertheless, mixture proportions in both data sets were dominated by Southeastern Alaska stocks. Nanopore assignments



**FIGURE 5** Principal coordinate analysis (PCoA) of SNP calls of individuals passing the threshold in all data sets. SNP calls based on Nanopore sequences generated during the IYS expedition are shown in blue (“nano\_IYS”), and the same individuals reanalysed upon return using the same workflow are 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”)

tended to overestimate the contribution to this stock as well as Lower Stikine River (LSTK) stocks. Many of the individuals assigned to these stocks using the Nanopore were assigned to the adjacent stocks of Lower Hecate Strait and Haro Strait (HecLow + HStr) as well as Southern Coastal Streams, Queen Charlotte Strait, Johnston Strait and Southern Fjords (SC + SFj) on the Ion Torrent platform (Figure 6, Table 1). Individuals from stocks well represented in the database such as the Columbia River were confidently assigned to the appropriate stock on both platforms. However, Z-scores calculated by rubias during stock assignment, which are an indirect measure of how well the SNP call matches individuals in the baseline data set of both, indicated that the Nanopore and the Ion Torrent data showed large deviations from the normal distribution, suggesting that many individuals assayed are not well represented in the database (Figure S5) (Moran & Anderson, 2019). Ion Torrent data show two peaks, one overlying the expected normal distribution and a second peak lying outside the normal distribution. This suggests that

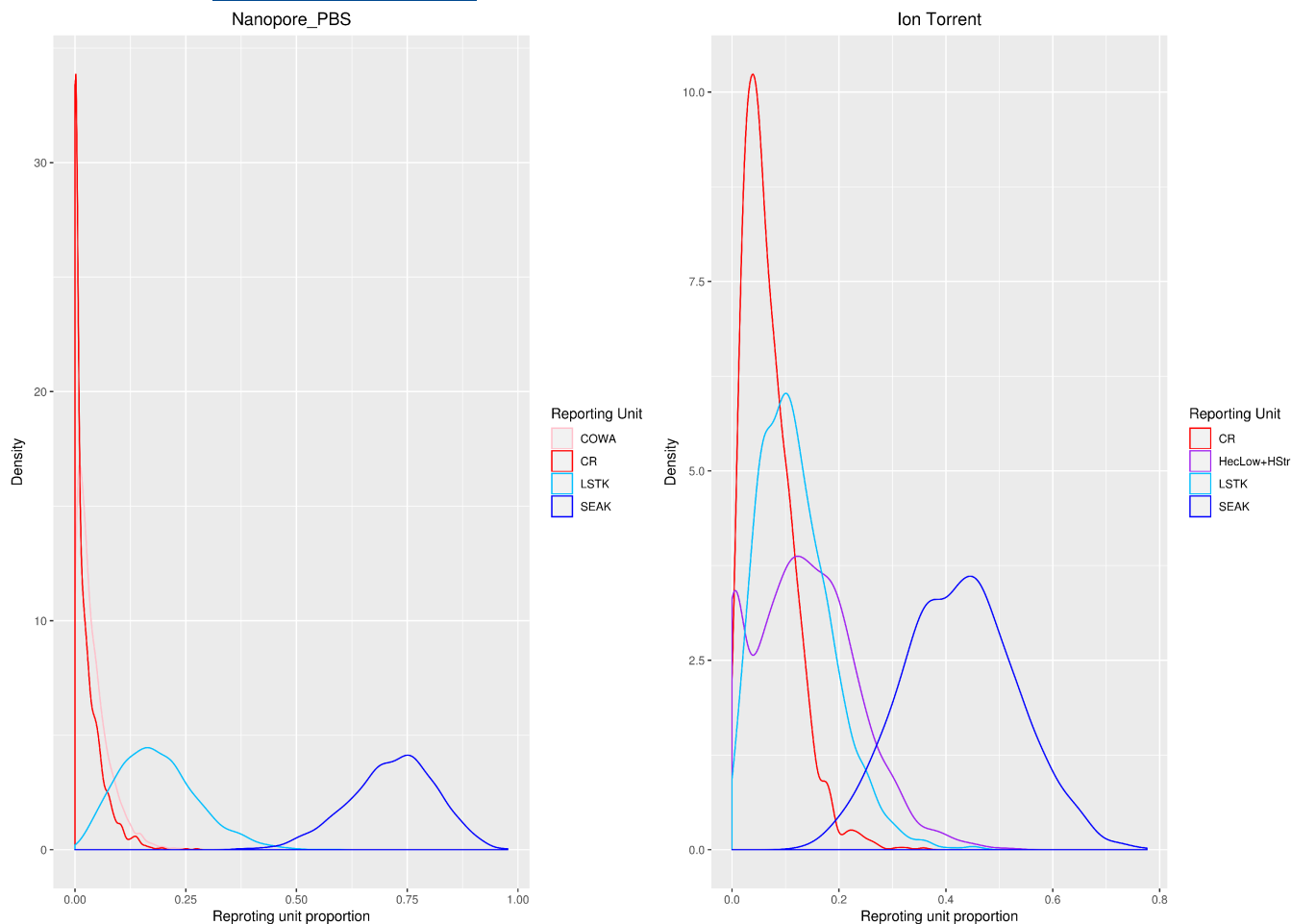
about half of the individuals were not from populations that are well represented in the database (Figure S5). Similarly, Nanopore-based assignments showed even more aberrant distributions, presumably due to the additive effects of the sequencing platform introducing bias on top of poor baseline representation (Figure S5). The poor database representation could cause small differences in SNP calls to cause alternative assignments.

## 4 | DISCUSSION

### 4.1 | Nanopore sequencing enables remote in-field SNP genetic stock identification

Here, we present the first proof-of-concept study demonstrating the feasibility of using the portable Oxford Nanopore min-ION sequencer for remote in-field genetic stock identification by





**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

**TABLE 1** Relative proportion of top reporting units (contribution >3%) to the overall mixture of coho salmon

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 + HStr	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			

Note: Only individuals that had a successful stock ID on all three GSI runs are included. Reporting units: SEAK: Southeast Alaska; LSTK: Lower Stikine River; NCS: North Coast Streams (BC); 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; LNASS: Lower Nass River; WVI: West Vancouver Island; OR: Oregon.

SNP sequencing of Pacific salmon. We developed a rapid sample processing workflow that relied on amplicon concatenation to increase throughput. With this workflow, we performed genetic

stock identification on 75 coho salmon onboard a research vessel in the Gulf of Alaska, with minimal equipment during two runs. Genetic stock identification of all 80 captured coho salmon in

a single run using the mobile platform resulted in stock assignment for 50 individuals at 67% concordance with state of the art laboratory-based pipelines.

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

The inherent low fidelity of the Nanopore platform using R9-type flow cells relative to other sequencing technologies, specifically around homopolymer tracts, proved to be the major shortcoming, limiting both the actual SNP calling accuracy, causing comparatively low repeatability, as well as the throughput, by necessitating a higher alignment coverage due to the high error rate (Cornelis et al., 2017). The low fidelity of the Nanopore sequences was specifically apparent when comparing it with the established sequencing platform for genetic stock identification by SNP sequencing, the Ion Torrent Proton sequencer (Beacham et al., 2017). The Ion Torrent short read sequencer routinely outperformed the Nanopore sequencer, both in accuracy and in throughput. The latter is a major restricting factor of the Nanopore platform due to a limited number of available sequencing pores inherent to the platform. While we compensated for this limitation by concatenating amplicons, to generate several amplicon sequences per Nanopore read, the efficiency of this approach was modest, yielding only a two-fold increase in throughput at present. Further, the needs for concatenation and higher inputs required several PCR amplification steps that could have contributed to the observed shifts in allele frequencies, leading to differing assignments on the different platforms. Turnaround time in the present study was mostly restricted by the computational capacity of the portable laptop used for the computational analysis. Specifically, base calling by translating the raw electrical signal recorded by the minION sequencer into fastq nucleotide reads proved to be the most time-consuming step, requiring up to several days in computing time.

However, despite the limitations associated with the Nanopore platform described above, the stock composition of coho salmon in the Gulf of Alaska also confounded accuracy and fidelity of stock assignment. Most importantly, most salmon sampled and assessed during the Gulf of Alaska expedition were assigned to Southeastern Alaska and adjacent British Columbia coast stocks (SEAK, HeLow + HStr, SC + SFj). These stocks are poorly represented in the queried baseline and stocks from northern Alaska are very sparse so that fish from such origin are often assigned to the SEAK with poor confidence. This meant that even on the Ion Torrent platform, assignment probabilities were low, causing small differences in SNP content between the two platforms that led to alternating assignment between these stocks (i.e., SEAK assignment on Nanaopore being assigned to HeLow + HStr and SC + SFj on Ion

Torrent). Indeed, stock assignment on the Ion Torrent platform using an updated and expanded baseline and primer set resulted in high-confidence assignment of many of these individuals to Kynoch and Mussel Inlets, a spatially close reporting unit on the Northern BC coast that was poorly represented in the original baseline (C. Neville, personal communication). This suggests that new SNP loci included in the updated primer set and baseline were able to resolve these stocks at higher confidence and assign them to the appropriate stock (Beacham et al., 2020). Fortunately, all of the current limitations mentioned above can be addressed in further development and we expect significant improvements in all fields, ultimately delivering a high-throughput, real-time, in-field sequencing platform.

#### 4.2 | Advances to the Nanopore platform, sample preparation, as well as computational infrastructure will improve turnaround, throughput and fidelity

While we were successful in providing a proof-of-principle study demonstrating that the Nanopore platform is capable of in-field genotyping, the throughput, fidelity and turnaround remained below the level needed to put this platform into standard operation for GSI by SNP genotyping. Several modifications in the workflow are planned to improve the throughput. Currently, barcoding relies on inefficient blunt-end ligation of the barcoding adapters to the PCR amplicons, leading to up to 50% unbarcoded amplicons and therefore wasting a large portion of sequencing capacity. Including the ligation adapter sequences needed to add the barcodes in the PCR primers will improve the efficacy of barcoding by circumventing the inefficient and laborious blunt-end ligation. This will improve sequencing throughput, while at the same time speeding up the sample preparation by approximately 1 hr. Next, concatenation efficiency is currently relatively low, increasing throughput only two-fold. While large concatemers approaching 10 kb were observed, they were relatively rare. Optimized concatenation conditions by adjusting the reaction conditions such as annealing temperature and duration should exponentially improve throughput by increasing both the relative abundance of concatenated amplicons and the total length of concatemers. Further workflow improvements could include pre-aliquoting of DNA extraction solution, barcodes and primers, as well as bead cleaning materials in 96-well plates before heading into the field, which should reduce an additional 2 h of sample preparation, as well as reduce the risk of cross-contamination in the field. Together, these improvements should bring the total sample preparation time to about 10 h, with approximately half the time being hands-on.

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

An additional issue for using Nanopore sequencing is the low accuracy of the sequencing platform at the time of this project using the R9 flow cells. This low accuracy requires excessively high alignment coverage at SNP locations to ensure accurate SNP calling. However, newer Nanopore flow cells promise greatly increased accuracy (e.g., 99.999% for R10) due to “a longer barrel and dual reader head” and have recently become available. This updated flow cell technology is therefore expected to greatly improve sequencing accuracy and possibly allow the lowering of alignment thresholds for SNP calling, thereby increasing the throughput more than two-fold. Improvements to the SNP calling pipeline might enable the identification and exclusion of erroneous SNP calls due to the ability to calculate the *p*-error associated with SNP calls, thereby increasing accuracy and repeatability. Finally, in selecting SNP loci for inclusion in GSI baselines, consideration of the types of sequences that are most problematic for Nanopore sequencing (e.g., homopolymer tracts) could go a long way to improving performance across platforms. Testing power in coastwide baselines once these problematic loci are excluded will be an important future step. Extrapolating the above-mentioned improvements would improve the current throughput of 96 individuals per flow cell by more than an order of magnitude, thereby enabling cost-effective real-time and/or field-based application of the platform.

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

## ACKNOWLEDGEMENTS

The authors would like to thank the following individuals for their contribution to the expedition and to the manuscript: Richard Beamish, Brian Riddell and the NPAFC secretariat for the organization of the 2019 Gulf of Alaska expedition; the entire scientific crew of the 2019 GoA expedition: Evgeny Pakhomov, Gerard Foley, Brian P. V. Hunt, Arkadii Ivanov, Hae Kun Jung, Gennady Kantakov, Anton Khleborodov, Chrys Neville, Vladimir Radchenko, Igor Shurpa, Alexander Slabinsky, Shigehiko Urawa, Anna Vazhova, Vishnu Suseelan, Charles Waters, Laurie Weitkamp and Mikhail

Zuev; the crew of the research vessel *Professor Kaganovskiy*; Charlie Waters for providing an R script for catch visualization; Chrys Neville for the contribution of catch data. This research was supported by the Pacific Salmon Commission, Pacific Salmon Foundation, and Fisheries and Oceans Canada and the Canadian Coast Guard (DFO CCG). C.M.D. was supported by a fellowship through the Pacific Salmon Foundation and MITACS.

## AUTHOR CONTRIBUTIONS

C.M.D., B.J.G.S. and K.M.M. designed research. C.M.D. performed research. T.J.M., C.W., K.J., K.L.F., E.B.R. and T.D.B. contributed new reagents or analytical tools. C.M.D., B.J.G.S. and E.B.R. analysed data. C.M.D., B.J.G.S. and K.M.M. wrote the paper.

## DATA AVAILABILITY STATEMENT

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

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

**Genetic Data:** All raw Nanopore sequence reads analysed in this paper are deposited in the SRA under BioProject: PRJNA796718 (SRR17593964–SRR17593966).

**Sample metadata:** Metadata on the individuals in this study are also stored associated with BioProject: PRJNA796718 under BioSamples SAMN24907542–SAMN24907622.

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

## BENEFIT SHARING STATEMENT

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

## ORCID

Christoph M. Deeg  <https://orcid.org/0000-0002-4459-9372>

Ben J. G. Sutherland  <https://orcid.org/0000-0002-2029-9893>

Terry D. Beacham  <https://orcid.org/0000-0003-0987-8445>

## REFERENCES

- Atlas, W. I., Ban, N. C., Moore, J. W., Tuohy, A. M., Greening, S., Reid, A. J., & Connors, K. (2021). Indigenous systems of management for culturally and ecologically resilient Pacific salmon (*Oncorhynchus* spp.) fisheries. *BioScience*, 71(2), 186–204.
- Beacham, T. D., Wallace, C. G., Jonsen, K., McIntosh, B., Candy, J. R., Rondeau, E. B., Moore, J.-S., Bernatchez, L., & Withler, R. E. (2020). Accurate estimation of conservation unit contribution to Coho Salmon mixed-stock fisheries in British Columbia, Canada using direct DNA sequencing for single nucleotide polymorphisms. *Canadian Journal of Fisheries and Aquatic Sciences*. <https://www.nrcresearchpress.com/doi/abs>

- Beacham, T. D., Wallace, C., MacConnachie, C., Jonsen, K., McIntosh, B., Candy, J. R., Devlin, R. H., & Withler, R. E. (2017). Population and individual identification of Coho Salmon in British Columbia through parentage-based tagging and genetic stock identification: An alternative to coded-wire tags. *Canadian Journal of Fisheries and Aquatic Sciences*, 74(9), 1391–1410. <https://doi.org/10.1139/cjfas-2016-0452>
- Beacham, T. D., Wallace, C., MacConnachie, C., Jonsen, K., McIntosh, B., Candy, J. R., & Withler, R. E. (2018). Population and individual identification of Chinook Salmon in British Columbia through parentage-based tagging and genetic stock identification with single nucleotide polymorphisms. *Canadian Journal of Fisheries and Aquatic Sciences*, 75(7), 1096–1105. <https://doi.org/10.1139/cjfas-2017-0168>
- Campbell, N. R., Harmon, S. A., & Narum, S. R. (2015). Genotyping-in-thousands by sequencing (GT-Seq): A cost effective SNP genotyping method based on custom amplicon sequencing. *Molecular Ecology Resources*, 15(4), 855–867. <https://doi.org/10.1111/1755-0998.12357>
- Cederholm, C. J., Kunze, M. D., Murota, T., & Sibatani, A. (1999). Pacific salmon carcasses: Essential contributions of nutrients and energy for aquatic and terrestrial ecosystems. *Fisheries*, 24(10), 6–15.
- Cook, R. C., & Guthrie, I. (1987). In-Season stock identification of Sockeye Salmon (*Oncorhynchus Nerka*) using scale pattern recognition. *Canadian Special Publication of Fisheries and Aquatic Sciences/Publication Speciale Canadienne Des Sciences Halieutiques Et Aquatiques*, 96, 327–334.
- Cornelis, S., Gansemans, Y., Deleye, L., Deforce, D., & Van Nieuwerburgh, F. (2017). Forensic SNP genotyping using Nanopore MinION sequencing. *Scientific Reports*, 7(February), 41759. <https://doi.org/10.1038/srep41759>
- Dann, T. H., Habicht, C., Baker, T. T., & Seeb, J. E. (2013). Exploiting genetic diversity to balance conservation and harvest of migratory salmon. *Canadian Special Publication of Fisheries and Aquatic Sciences/Publication Speciale Canadienne Des Sciences Halieutiques Et Aquatiques*, 70(5), 785–793. <https://doi.org/10.1139/cjfas-2012-0449>
- Gilbey, J., Wennevik, V., Bradbury, I. R., Fiske, P., Hansen, L. P., Jacobsen, J. A., & Potter, T. (2017). Genetic stock identification of Atlantic salmon caught in the Faroese fishery. *Fisheries Research*, 187(March), 110–119. <https://doi.org/10.1016/j.fishres.2016.11.020>
- Hinch, S. G., Cooke, S. J., Farrell, A. P., Miller, K. M., Lapointe, M., & Patterson, D. A. (2012). Dead fish swimming: a review of research on the early migration and high premature mortality in adult Fraser river sockeye salmon *Oncorhynchus Nerka*. *Journal of Fish Biology*, 81(2), 576–599. <https://doi.org/10.1111/j.1095-8649.2012.03360.x>
- Jefferts, K. B., Bergman, P. K., & Fiscus, H. F. (1963). A coded wire identification system for macro-organisms. *Nature*, 198(4879), 460–462. <https://doi.org/10.1038/198460a0>
- Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics*, 25(14), 1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., & Durbin, R. (2009). The Sequence Alignment/Map Format and SAMtools. *Bioinformatics*, 25(16), 2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>
- Lichtatowich, J. (2001). *Salmon without rivers: A history of the pacific salmon crisis*. Island Press.
- Mikheyev, A. S., & Tin, M. M. Y. (2014). A First look at the Oxford Nanopore MinION sequencer. *Molecular Ecology Resources*, 14(6), 1097–1102. <https://doi.org/10.1111/1755-0998.12324>
- Miller, K. M., Teffer, A., Tucker, S., Li, S., Schulze, A. D., Trudel, M., Juanes, F., Tabata, A., Kaukinen, K. H., Ginther, N. G., Ming, T. J., Cooke, S. J., Hipfner, J. M., Patterson, D. A., & Hinch, S. G. (2014). Infectious disease, shifting climates, and opportunistic predators: Evolutionary factors potentially impacting wild salmon declines. *Evolutionary Applications*, 7(7), 812–855. <https://doi.org/10.1111/eva.12164>
- Miller, K. M., Withler, R. E., & Beacham, T. D. (1996). Stock identification of Coho Salmon (*Oncorhynchus Kisutch*) using minisatellite DNA variation. *Canadian Journal of Fisheries and Aquatic Sciences*, 53(1), 181–195.
- Moran, B. M., & Anderson, E. C. (2019). Bayesian Inference from the Conditional Genetic Stock Identification Model. *Canadian Journal of Fisheries and Aquatic Sciences*, 76(4), 551–560. <https://doi.org/10.1139/cjfas-2018-0016>
- Ozerov, M., Vasemägi, A., Wennevik, V., Diaz-Fernandez, R., Kent, M., Gilbey, J., Prusov, S., Niemelä, E., & Vähä, J.-P. (2013). Finding markers that make a difference: DNA Pooling and SNP-arrays identify population informative markers for genetic stock identification. *PLoS One*, 8(12), e82434. <https://doi.org/10.1371/journal.pone.0082434>
- Paradis, E., & Schliep, K. (2019). Ape 5.0: An environment for modern phylogenetics and evolutionary analyses in R. *Bioinformatics*, 35(3), 526–528. <https://doi.org/10.1093/bioinformatics/bty633>
- Quick, J., Loman, N. J., Duraffour, S., Simpson, J. T., Severi, E., Cowley, L., Bore, J. A., Koundouno, R., Dudas, G., Mikhail, A., Ouédraogo, N., Afrough, B., Bah, A., Baum, J. H. J., Becker-Ziaja, B., Boettcher, J. P., Cabeza-Cabrerizo, M., Camino-Sánchez, Á., Carter, L. L., ... Carroll, M. W. (2016). Real-time, portable genome sequencing for ebola surveillance. *Nature*, 530(7589), 228–232. <https://doi.org/10.1038/nature16996>
- Robinson, J. T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E. S., Getz, G., & Mesirov, J. P. (2011). Integrative genomics viewer. *Nature Biotechnology*, 29(1), 24–26. <https://doi.org/10.1038/nbt.1754>
- Schlecht, U., Mok, J., Dallett, C., & Berka, J. (2017). ConcatSeq: A method for increasing throughput of single molecule sequencing by concatenating Short DNA fragments. *Scientific Reports*, 7(1), 5252. <https://doi.org/10.1038/s41598-017-05503-w>
- Winans, G. A., Aebersold, P. B., Urawa, S., & Varnavskaya, N. V. (1994). Determining continent of origin of chum salmon (*Oncorhynchus keta*) using genetic stock identification techniques: status of allozyme baseline in Asia. *Canadian Journal of Fisheries and Aquatic Sciences*, 51(S1), 95–113.
- Wood, C. C., Rutherford, D. T., & McKinnell, S. (1989). Identification of Sockeye Salmon (*Oncorhynchus Nerka*) stocks in mixed-stock fisheries in British Columbia and southeast Alaska using biological markers. *Canadian Journal of Fisheries and Aquatic Sciences*, 46(12), 2108–2120.

## SUPPORTING INFORMATION

Additional supporting information may be found in the online version of the article at the publisher's website.

**How to cite this article:** Deeg, C. M., Sutherland, B. J. G., Ming, T. J., Wallace, C., Jonsen, K., Flynn, K. L., Rondeau, E. B., Beacham, T. D., & Miller, K. M. (2022). In-field genetic stock identification of overwintering coho salmon in the Gulf of Alaska: Evaluation of Nanopore sequencing for remote real-time deployment. *Molecular Ecology Resources*, 22, 1824–1835. <https://doi.org/10.1111/1755-0998.13595>