

Article

# MinION-Based DNA Barcoding of Preserved and Non-Invasively Collected Wildlife Samples

Adeline Seah <sup>1,†</sup> , Marisa C.W. Lim <sup>1,\*,†</sup> , Denise McAloose <sup>1</sup>, Stefan Prost <sup>2,3</sup>  and Tracie A. Seimon <sup>1</sup> 

<sup>1</sup> Zoological Health Program, Wildlife Conservation Society, Bronx Zoo, 2300 Southern Blvd, Bronx, NY 10460, USA; biodiversityconnections@gmail.com (A.S.); dmcaloose@wcs.org (D.M.); tseimon@wcs.org (T.A.S.)

<sup>2</sup> LOEWE-Centre for Translational Biodiversity Genomics, Senckenberg Nature Research Society, 60325 Frankfurt, Germany; stefanprost.research@protonmail.com

<sup>3</sup> South African National Biodiversity Institute, National Zoological Garden, Pretoria 0001, South Africa

\* Correspondence: marisa\_lim@pacbell.net

† These authors contributed equally.

Received: 27 March 2020; Accepted: 16 April 2020; Published: 18 April 2020



**Abstract:** The ability to sequence a variety of wildlife samples with portable, field-friendly equipment will have significant impacts on wildlife conservation and health applications. However, the only currently available field-friendly DNA sequencer, the MinION by Oxford Nanopore Technologies, has a high error rate compared to standard laboratory-based sequencing platforms and has not been systematically validated for DNA barcoding accuracy for preserved and non-invasively collected tissue samples. We tested whether various wildlife sample types, field-friendly methods, and our clustering-based bioinformatics pipeline, SAIGA, can be used to generate consistent and accurate consensus sequences for species identification. Here, we systematically evaluate variation in cytochrome b sequences amplified from scat, hair, feather, fresh frozen liver, and formalin-fixed paraffin-embedded (FFPE) liver. Each sample was processed by three DNA extraction protocols. For all sample types tested, the MinION consensus sequences matched the Sanger references with 99.29%–100% sequence similarity, even for samples that were difficult to amplify, such as scat and FFPE tissue extracted with Chelex resin. Sequencing errors occurred primarily in homopolymer regions, as identified in previous MinION studies. We demonstrate that it is possible to generate accurate DNA barcode sequences from preserved and non-invasively collected wildlife samples using portable MinION sequencing, creating more opportunities to apply portable sequencing technology for species identification.

**Keywords:** Biomeme; Chelex; DNA barcoding; FFPE; field-friendly; MinION; non-invasive sampling

## 1. Introduction

Wildlife health and conservation initiatives benefit tremendously from genetic methods of species identification for infectious disease screening [1,2], detecting illegally traded wildlife products [3], uncovering food label fraud [3–5], and documenting understudied biodiversity [6]. One major challenge for wildlife molecular studies is obtaining fresh samples from live or dead wild animals. Such endeavors can be logistically challenging, generally involving highly skilled teams, detailed planning, and acquisition of permissions from local, regional, and international partners and governmental agencies for animal handling, sample collection, and sample transfer for molecular testing. Consequently, environmental samples [7,8] and animal samples that can be collected non-invasively (e.g. hair, feathers, scat, etc.) [9–11] are increasingly being used for ecological studies, wildlife health assessments, and characterizing biodiversity. Non-invasively collected samples are easier to obtain

than fresh organ tissues, but may contain PCR inhibitors, have lower DNA yields, or be degraded from environmental exposure [10,12–14]. Archived historical wildlife samples, often preserved in formalin, also offer a unique opportunity to obtain genetic information [15]. However, challenges for molecular studies include formalin-related fragmentation and DNA cross-linking [16,17].

DNA barcoding is a common molecular technique for species identification [18,19] that had been traditionally carried out via standard lab-based sequencing equipment until recent developments in portable technology. The Oxford Nanopore Technologies (ONT) MinION sequencer is currently the only available portable sequencer. Although nanopore sequencing is known to have higher raw sequence error rates in comparison to standard short read sequencing platforms such as Illumina or BGI-Seq, particularly at homopolymeric regions [20,21], significant improvements in the accuracy of MinION sequencing chemistry has led to its recent rise in popularity for field applications (reviewed in [22]). This sequencer is especially useful in situations where there is a lack of access to sequencing facilities or when sample export is difficult. The MinION also has a lower investment cost and shorter turnaround times than traditional sequencing platforms (e.g., Sanger, Illumina).

MinION DNA barcoding studies have primarily used laboratory-based QIAGEN®kits for reliable and pure DNA extraction products (e.g., [23–25]). To expand the potential for portable sequencing applications, field-friendly DNA extraction methods can be used to reduce lab equipment requirements. While field-friendly DNA extraction methods are often less effective at producing DNA of high concentration and purity levels, MinION DNA barcoding has been successfully performed using QuickExtract™ solution (Lucigen, Middleton, USA), which only requires a heat source [26]. The Chelex®100 resin (Bio-Rad Inc., Hercules, USA) extraction method similarly only requires a heat source (e.g., heat block or PCR thermocycler), but is less expensive and has not been tested for MinION sequencing so far. Both methods have short protocols, but do not remove cellular debris or PCR inhibitors, which can affect downstream applications [27,28]. The Biomeme M1 Sample Prep™ Kit (Biomeme Inc., Philadelphia, USA) is another DNA extraction kit developed for field use. While more expensive than either QuickExtract or Chelex methods, the Biomeme kit includes all necessary components and both protein and salt wash steps to remove impurities. Studies have shown that Biomeme-extracted samples have higher levels of inhibitors compared to Qiagen extractions, and thus require additional dilution steps [8,29].

To date, MinION DNA barcoding pipelines have used either *de novo* assembly [23,24], clustering-based [25], or alignment [26,30] methods to generate consensus sequences for species identification. Assembly approaches generally work more consistently for longer barcodes ( $\approx 1$  kb), as the underlying software were originally designed for assembling long reads for genome assemblies rather than amplicons. To date, published clustering [25] or alignment [26,30] pipelines use subsets of the data (100–200 reads) to generate scaffolds for read error correction. While these approaches may work for high quality sequence data, such data subsets could include more sequence error bias in lower quality datasets. Thus, we developed a clustering-based pipeline, SAIGA (<https://github.com/marisalim/Saiga>), with software specifically designed for error prone MinION reads that processes data regardless of barcode length and has no limits on the number of reads that can be clustered, thus maximizing the use of demultiplexed reads for downstream species identification analysis.

In this study, we systematically evaluate the accuracy of the MinION for DNA barcoding across a range of wildlife sample types, including two field-friendly DNA extraction approaches. We sequenced a short fragment of the commonly used mitochondrial cytochrome b (Cytb) gene from scat, hair, feather, fresh frozen liver, and formalin-fixed paraffin embedded (FFPE) liver. For each sample type, we compared the accuracy of Cytb consensus sequences for three different DNA extraction methods: QIAGEN silica membrane-based kits (Qiagen Inc., Germantown, USA), Chelex 100 resin (Bio-Rad Inc., Hercules, USA), and the Biomeme M1 Sample Prep Kit (Biomeme Inc., Philadelphia, USA). All analyses were conducted with SAIGA. We demonstrate that MinION sequencing can be used with field-friendly extraction methods to accurately identify wildlife species from a variety of sample types.

## 2. Materials and Methods

### 2.1. Sample Collection

For this study, scat, hair, feather, fresh frozen liver, and FFPE liver samples were collected opportunistically during necropsy examinations from a snow leopard (*Panthera uncia*) and a cinnamon teal (*Anas cyanoptera*) from a zoological collection. The FFPE liver samples were part of a suite of tissues that were collected, stored in 10% neutral buffered formalin, and subsequently processed and paraffin-embedded for histologic examination and routine tissue archiving. Fresh liver, scat, hair, and feather samples were frozen (−80 °C) immediately after collection.

### 2.2. DNA Extraction

DNA was extracted from each sample type using three different approaches: (1) Qiagen (QIAamp®DNA minikit or QIAamp®DNA Stool Mini Kit, Qiagen Inc., Germantown, MD, USA); (2) Chelex 100 Resin (Bio-Rad Inc., Hercules, CA, USA); and (3) Biomeme M1 Sample Prep Kit for DNA (Biomeme Inc., Philadelphia, PA, USA). DNA quantification is inaccurate for Chelex extracts due to the presence of cellular components, thus Chelex extracts were not quantified. All Qiagen and Biomeme extracts were quantified using the Qubit™ dsDNA High Sensitivity Kit on the Qubit™ 4 Fluorometer (Thermo Fisher Scientific, Waltham, MA, USA). The Qiagen, Chelex, and Biomeme extraction protocols are summarized for each tissue type in Text S1. All Qiagen, Biomeme DNA extracts with >10 ng/μL, and all Chelex extracts were run on a 1.0% gel to assess DNA fragmentation by sample type.

### 2.3. PCR and Library Preparation

#### 2.3.1. DNA Barcoding PCR—Round 1

Approximately 460 bp of the mitochondrial *Cytb* gene was amplified using primers mcb398 and mcb869 [31], with universal tailed sequences on each primer that are compatible with the ONT PCR Barcoding Expansion kit EXP-PBC001 (ONT, Oxford, UK) (Table S1). These primers were designed from an alignment of 67 animal species, and validated for mammals, reptiles, and birds [31].

PCR was carried out with 6.25 μL DreamTaq HotStart PCR Master Mix (Thermo Fisher, Waltham, MA, USA), 1.25 μL DNA template, and 2 μL of each primer (10 μM stock) in a final volume of 12.5 μL. Cycling conditions were: 95 °C for 3 min; 35 cycles of 95 °C for 30 s, 55 °C for 30 s, and 72 °C for 30 s; and a final extension of 72 °C for 5 min. All Chelex extractions were diluted for the DNA Barcoding PCR as described in Text S1. PCR products were purified using 1.8× Agencourt AMPure XP beads (Beckman Coulter, Indianapolis, IN, USA), tested for purity using the NanoDrop™ One spectrophotometer (Thermo Fisher Scientific, Waltham, MA, USA), and quantified fluorometrically using the Qubit dsDNA High sensitivity kit.

#### 2.3.2. Indexing PCR—Round 2

To attach dual ONT PCR index sequences to the *Cytb* amplicons, a second round of PCR was carried out with the ONT PCR Barcoding Expansion kit EXP-PBC001 for each sample with 25 μL KAPA Biosystems HiFi HotStart ReadyMix (2×) (Thermo Fisher Scientific, Waltham, USA), containing 25 ng of first-round PCR amplicon and 1 μL ONT PCR Barcode in a final volume of 50 μL. Cycling conditions were: 95 °C for 3 min; 11 cycles of 95 °C for 15 s, 62 °C for 15 s, and 72 °C for 15 s; and a final extension of 72 °C for 1 min. Hereafter, we refer to ONT PCR barcodes as ‘indexes’ to reduce confusion with the *Cytb* barcode. Indexed PCR products from round 2 were purified and tested for purity and quantity like round 1 products.

#### 2.3.3. Library Preparation

Samples were grouped into four libraries by sample type (FFPE, scat, hair/feather, frozen liver). For each library, purified indexed amplicons were pooled in equal ratios to produce 1.0–1.2 μg in a total

of 45  $\mu\text{L}$  nuclease-free water. Pooled libraries were next prepared using the ONT Ligation Sequencing kit SQK-LSK109 (ONT, Oxford, UK) with modifications to the manufacturer's instructions: 25  $\mu\text{L}$  of the pooled library was mixed with 3.5  $\mu\text{L}$  NEBNext Ultra II End-Prep Reaction buffer and 1.5  $\mu\text{L}$  Ultra II End-prep Enzyme mix (New England Biolabs, Ipswich, MA, USA), incubated for 10 min at room temperature, then 10 min at 65  $^{\circ}\text{C}$ . For adapter ligation, 15  $\mu\text{L}$  of the end-prepped library (not bead-purified) was mixed with 25  $\mu\text{L}$  Blunt/TA Ligase and 10  $\mu\text{L}$  Adapter Mix (AMX), incubated at room temperature for 20 min and eluted in a final volume of 12  $\mu\text{L}$  of Elution Buffer.

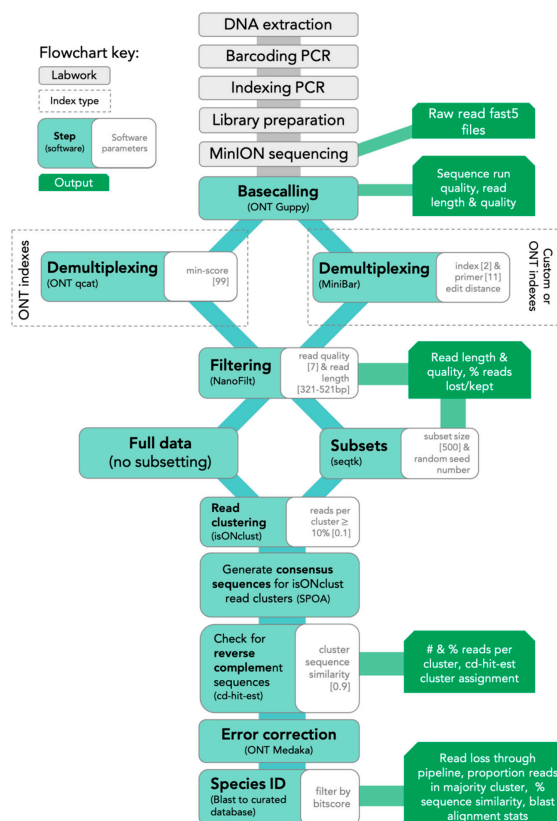
#### 2.4. Sequencing

The four libraries were split between two FLO-MIN106D R9.4.1 chemistry flow cells (ONT, Oxford, UK) to minimize bleed-through between experiments, with two libraries run on each flow cell—FAL19910: (1) FFPE (four samples, BC01-04), (2) scat (six samples, BC05-10); FAL19272: (1) hair/feather (six samples, BC01-06), (2) frozen liver (six samples, BC07-12). Flow cells were washed with Wash Solution A followed by the addition of Storage Buffer S according to the manufacturer's protocols. All libraries were sequenced for approximately 1 h to obtain at least 100,000 raw reads per sample.

For comparison to MinION sequences, Sanger sequencing in the forward and reverse directions was performed on all purified indexed amplicons (Eton Bioscience Inc. Newark, NJ, USA). Sanger consensus sequences were generated using Geneious Prime v2019.0.4 software (Biomatters LDT, Auckland, New Zealand).

#### 2.5. Bioinformatics

The SAIGA bioinformatics pipeline is available on GitHub (<https://github.com/marisalim/Saiga>) and steps are outlined in Figure 1. MinKNOW (ONT) was used for sequencing and the raw sequence data were basecalled using Guppy v3.5.1 (ONT) with basecalling model "dna\_r9.4.1\_450bps\_fast.cfg".



**Figure 1.** Lab and SAIGA bioinformatics pipeline flowchart. Bioinformatics software and parameters are indicated at each step.

### 2.5.1. Demultiplexing and Filtering

Assigning sequencing reads to the correct sample is a critical step to avoid mixing sample sequences within or between sequencing runs. Thus, we compared results from two demultiplexing programs: (1) qcat v1.1.0 (ONT, <https://github.com/nanoporetech/qcat>) and (2) MiniBar v0.21 [24]. The qcat software was built specifically for demultiplexing reads indexed with ONT's barcode kits, while MiniBar is a general demultiplexing software that allows any set of user-specified index and primer sequences. We used stringent demultiplexing filters based on software recommendations, sensitivity analyses, and to minimize incorrect read assignments. Qcat uses the epi2me demultiplexing algorithm and we trimmed adapter and index sequences with the trim option. Using the min-score option, demultiplexed reads with alignment scores <99 were removed prior to downstream analysis, where a score of 100 means every nucleotide of the index is correct. Lower min-score thresholds (i.e., 60–90) reduced downstream consensus sequence quality. In MiniBar, up to two nucleotide differences between reads were allowed for the index sequences and 11 nucleotide differences between primer sequences per software recommendations; MiniBar primarily uses the index sequence information to demultiplex and trim dual index and primer sequence.

After demultiplexing, reads were removed if they had mean Phred quality scores <7 and were longer or shorter than the target amplicon length ( $\approx$ 421 bp excluding primers) with a 100 bp buffer (321–521 bp) in NanoFilt v2.5.0 [32]. Following each of the above steps, we calculated and visualized read quality statistics for raw, demultiplexed, and filtered reads with NanoPlot v1.21.0 [32]. To standardize dataset size across the four sequencing experiments and to investigate the effect of read depth, we generated 100, 500, and 5000 random read subsets for each sample from the filtered demultiplexed read files. Hereafter, we refer to these subsets as 100R, 500R, and 5KR, respectively.

### 2.5.2. Read Clustering and Consensus Sequence Generation

To generate the consensus sequence for each sample, all reads were first clustered using isONclust v0.0.4 [33]. We chose isONclust over clustering tools previously used in nanopore-based DNA barcoding pipelines, such as VSEARCH (implemented in ONTrack, [25]), as it was specifically designed to work with error-prone long-read data and thus should be less affected by read errors and more efficient in cluster formation. Next, SAIGA outputs the number of reads per cluster, only retaining clusters with >10% of the total reads (user-defined). We implemented this step to minimize the inclusion of reads with high sequence error and possible contaminant reads in downstream analysis. Intermediate consensus sequences are then generated using SPOA v3.0.1 (<https://github.com/rvaser/spoa>), which is based on a partial order alignment (POA) algorithm [34]. SPOA also conducts error corrections, resulting in more accurate consensus sequences. The SPOA consensus sequences are then clustered using cd-hit-est v4.8.1 with a stringent similarity cutoff (0.9; user-defined) [35,36]. Since isONclust separates reads in different strand orientations, this second round of clustering groups reverse-complement SPOA consensus sequences, ensuring that more filtered reads are used for generating the final consensus sequence. The reads contributing to all SPOA consensus sequences that group with the majority isONclust cluster's SPOA consensus sequence are combined into a single file for mapping. SAIGA then maps these reads to the SPOA consensus sequence of the majority isONclust cluster for consensus polishing with ONT's Medaka software v0.10.0 (<https://github.com/nanoporetech/medaka>).

## 2.6. Consensus Accuracy and Analysis

The MinION consensus sequences were compared to Sanger sequences from the same sample using a nucleotide Blast search v2.8.1+ [37]. To assess and compare species identification results across tissue types, extraction methods, demultiplexing programs, and data subsets, the following were evaluated: (1) the percent of matching nucleotides between consensus and Sanger sequences, (2) the number of matching nucleotides between consensus and Sanger sequences, and (3) the proportion of filtered reads in the cluster used to generate final consensus sequence. Accurate species identification

was defined as those with >99% sequence similarity to the Sanger sequence and  $\approx 421$  bp of matching nucleotides. The proportion of demultiplexed reads contributing to the final consensus indicates how much data was used for species identification. For samples with consensus sequences generated from fewer than  $\approx 75\%$  of reads, we investigated the non-majority isONclust clusters for potential sequence error or contaminant reads. Finally, all MinION consensus and Sanger sequences across tissue types, extraction methods, demultiplexing software, and data subsets were aligned with Mafft v1.3.7 in Geneious Prime v2019.0.4 to identify common regions with sequence errors.

### 2.7. Data Availability

A representative Sanger sequence for both species is available on GenBank (MN823069-70), and MinION fastq files (basecalled, demultiplexed, and filtered) are available on NCBI Short Read Archive (BioProject: PRJNA594927, accessions: SRR10678113-SRR10678156). Raw MinION sequence data is available on the EBI European Nucleotide Archive (ERP119594).

## 3. Results

### 3.1. DNA Barcoding and Indexing PCR Performance

DNA concentrations were higher for Qiagen (0.8 to 59 ng/ $\mu$ L,  $n = 8$ ) compared to Biomeme (0.07 to 13.9 ng/ $\mu$ L,  $n = 8$ ) extractions (Table S2); Chelex samples were not quantified ( $n = 8$ ). Gel electrophoresis of Qiagen-extracted tissues show frozen liver and scat samples had high molecular weight genomic DNA, while FFPE samples were fragmented; hair and feather extracts were too faint to detect reliably (Figure S1). We were unable to detect high molecular weight nucleic acid in the Biomeme and Chelex-extracted samples (Figure S2). Despite variation in starting DNA concentration and the presence of low molecular weight fragments in some samples, we successfully barcoded and indexed 22 of 24 samples. The two samples that failed to amplify at the Barcoding PCR (Round 1) step were the snow leopard FFPE samples extracted by the Chelex and Biomeme protocols. The DNA concentration of DNA Barcoding PCR (Round 1) products after bead clean-up was <13.9 ng/ $\mu$ L with an average of 3.49 ng/ $\mu$ L. At these low DNA concentrations, NanoDrop purity of Barcoding Round 1 amplicons is highly variable and not reliable.

Two samples had less than 25 ng for Indexing PCR (Round 2). After bead clean-up, the concentration of the snow leopard liver/Chelex DNA Barcoding PCR (Round 1) product was much lower than expected (4.4 ng), despite having a bright agarose gel band. Nevertheless, this was sufficient for amplification in the Indexing PCR step. *Cytb* was also difficult to amplify from the snow leopard scat/Chelex, so amplicons from two DNA Barcoding (Round 1) PCR reactions were pooled for a total of 16 ng to proceed with Indexing PCR (Round 2). After the Indexing PCR (Round 2) bead clean-up, DNA concentrations were >19 ng/ $\mu$ L with an average of 80.92 ng/ $\mu$ L for all but the snow leopard liver/Chelex sample, which had 6.58 ng/ $\mu$ L. Average  $A_{260}/A_{280}$  ratios (1.82) and  $A_{260}/A_{280}$  ratios (1.96) indicated relatively pure samples for library preparation.

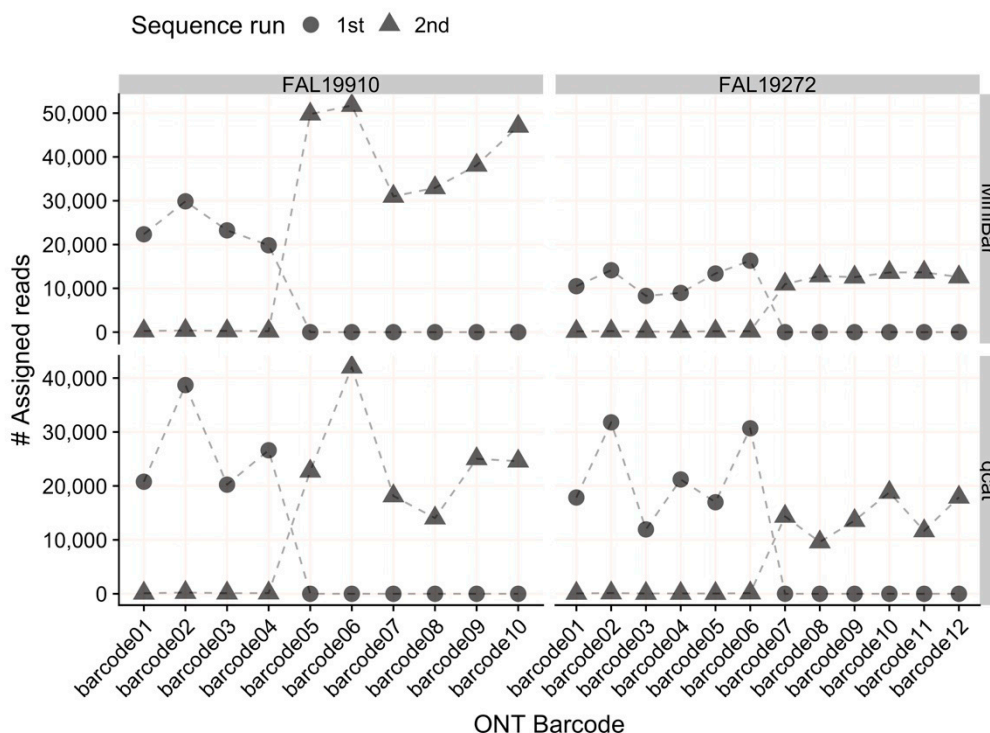
### 3.2. MinION and Sanger Sequencing Performance

Sequencing efficiency, also called pore occupancy, ranged from 72% to 80% and was evenly spread across flow cells for all MinION sequencing runs (Figure S3). We sequenced an average of  $\approx 752,856$  raw reads per run, with an average read length of  $\approx 597$  bp and read quality Phred score of 10.5 (Table S3, Figure S4).

We obtained clean Sanger sequences for 21 of 22 samples, all of which were 421 bp after primer trimming (Table S4). For all 21 samples, the Sanger sequences for each species were identical, regardless of tissue type or extraction method. We were unable to get a clean Sanger sequence for the snow leopard scat/Chelex sample. Therefore, we compared the MinION scat/Chelex consensus to the Sanger sequences from the other snow leopard samples for species identity.

### 3.3. Sequence Read Retention After Demultiplexing and Filtering

The average read quality and read lengths were similar across all samples demultiplexed with MiniBar or qcat (Tables S3 and S4). For all sequencing runs, both MiniBar and qcat correctly assigned demultiplexed reads only to the ONT indexes used in the Indexing PCR for each run (Figure 2). Due to the stringent demultiplexing thresholds, the majority of read data loss occurred during the demultiplexing step (84.07% reads lost on average; Table S3). After read quality and length filtering, we retained nearly all demultiplexed reads (95.6% reads retained on average; Figure S5, Table S3). On average, samples had more than 20,000 demultiplexed and filtered reads for downstream analyses (Table S4). In general, MiniBar-demultiplexed datasets retained more reads than qcat-demultiplexed datasets after filtering (Figure S5). The only sample that retained fewer than 90% of reads after filtering was the cinnamon teal scat/Biomeme sample demultiplexed with MiniBar (68.90% reads retained).



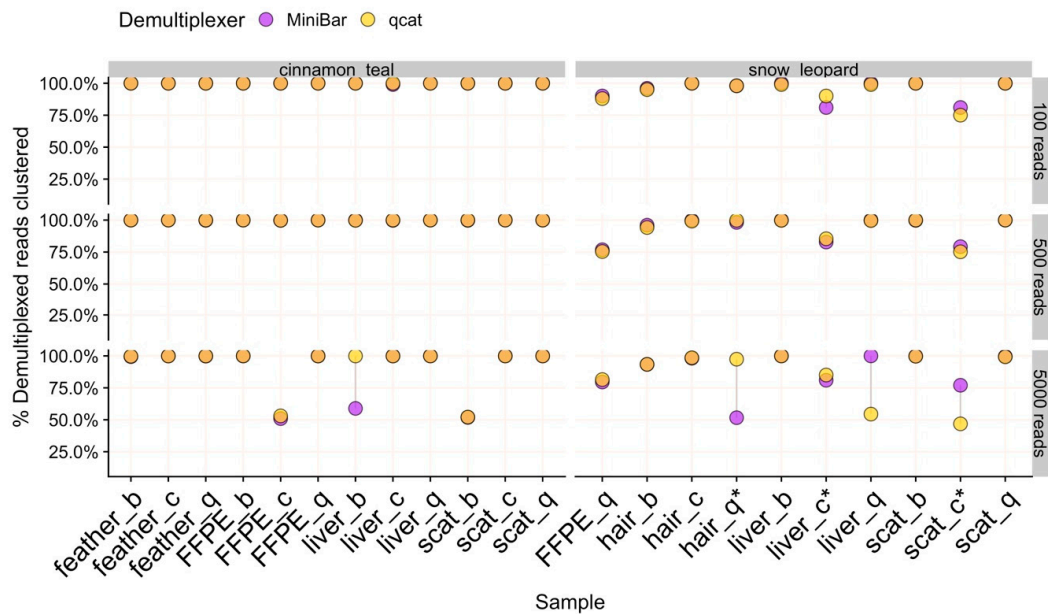
**Figure 2.** The number of reads assigned to each Oxford Nanopore Technologies (ONT) index (01–12) per flow cell by MiniBar and by qcat. For flow cell FAL19910, the 1st sequencing run used indexes 01–04 and the 2nd run used indexes 05–10. For flow cell FAL19272, the 1st sequence run used indexes 01–06 and the 2nd run used indexes 07–12.

### 3.4. Read Clustering Proportions and Cluster Species Identity

For nearly all data subsets, isONclust grouped reads into two clusters, one each for forward and reverse-complement oriented reads. In these cases, 100% of filtered reads formed a single cluster after cd-hit clustering to merge potential reverse-complements and thus, all reads were used to produce the consensus sequence for final species identification (Figure 3).

However, the remaining 18 data subsets (of 132) had more than two isONclust clusters, which fall into two categories: (1) samples where fewer than 60% of reads were used for final consensus generation due to sequence error and (2) samples with clusters containing contaminant reads (Table S5). In 5KR subsets for three cinnamon teal (FFPE/Chelex, liver/Biomeme, scat/Biomeme) and two snow leopard (hair/Qiagen, liver/Qiagen) samples, the second largest isONclust cluster contained reads that best match the same species as the majority cluster. While SPOA consensus sequences for these two clusters remained separate after cd-hit-est clustering, likely due to sequencing error (Table S5), species identification was successful for these five 5KR subset samples using only  $\approx 50\%$  of the reads to build

the consensus. In comparison, 100% of the reads clustered for the 100R and 500R subsets for these samples, suggesting that the 5KR subsample contained slightly more variation in read quality than the smaller subsets.



**Figure 3.** The percent of demultiplexed reads used to generate the final consensus sequence for 100R, 500R, and 5KR subsets for each species. Samples are labeled by tissue type and extraction method (b = biomeme, c = chelex, q = qiagen). Points are linked by a grey line to show difference in values from demultiplexers. Overlapping areas in orange indicate similar results for Minibar and qcat analyses. Asterisks (\*) indicate samples with cinnamon teal contamination.

We detected low to medium levels of cinnamon teal reads in three snow leopard samples: hair/Qiagen, scat/Chelex, and liver/Chelex, where the full set of demultiplexed reads contained 3.9%, 22.0%, and 14.4% teal reads, respectively. There were no teal contaminant reads, and hence no teal read clusters, in the snow leopard hair/Qiagen sample for all subsets. In contrast, the proportions of reads used to generate final consensus for all subsets of the snow leopard scat/Chelex and liver/Chelex samples were reduced to 75–85% of reads (Table S5). Recovery of DNA Barcoding PCR (Round 1) products was low for these two samples. However, our pipeline’s filtering and clustering procedures were able to correctly identify these samples as snow leopard because reads with high sequence errors and contaminant reads were not included in downstream analysis. There were no cinnamon teal reads in the rest of the snow leopard samples, and no snow leopard reads in any cinnamon teal samples.

### 3.5. Consensus Sequence Generation

The average proportion of reads used and consensus sequence lengths were comparable between sample types, extraction methods, subsets, and demultiplexers (Table 1 and Table S6). In general, SAIGA retained similar proportions of reads to generate consensus sequences across samples extracted by the Biomeme and Chelex methods as compared to the gold standard Qiagen-extracted samples (Figure 3, Table 1 and Table S6). In two cases, greater proportions of reads were used for the snow leopard liver and hair samples extracted with the Biomeme and Chelex protocols compared to the Qiagen-extract of the same tissue type. For samples where the consensus sequence length differed by demultiplexer, MiniBar subsets produced slightly longer sequences than qcat subsets (Figure S6).

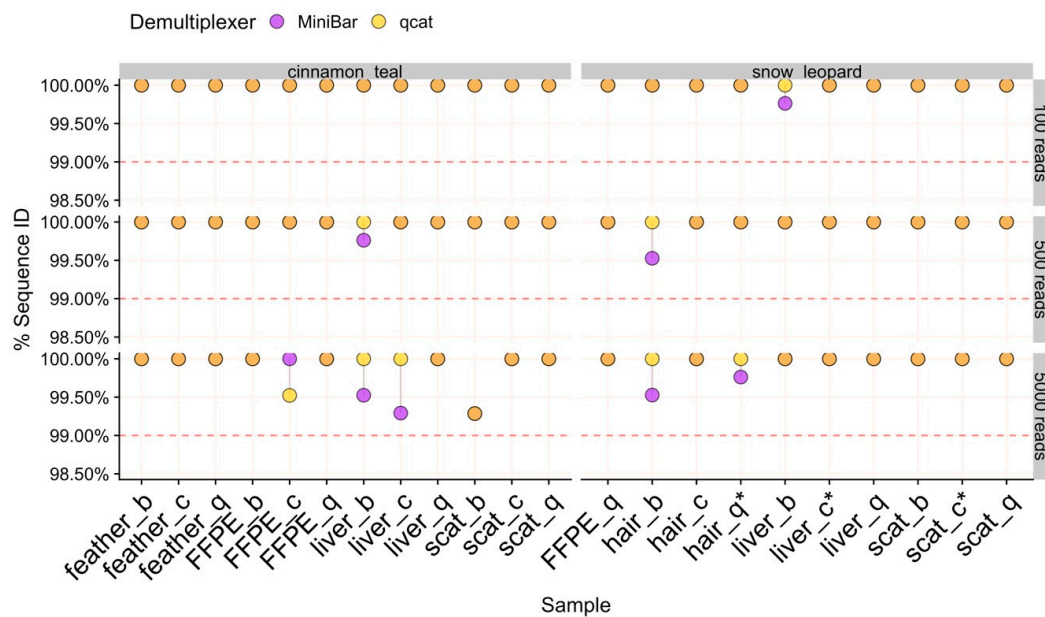


**Table 1.** Average and standard deviation (sd) for percent sequence similarity to Sanger sequence, length of matching nucleotides, and number and percent of demultiplexed reads used for the final consensus sequence from 100R, 500R, or 5KR read subsets demultiplexed with MiniBar or qcat. Statistics were calculated across all tissue types and extraction method samples.

Subset	Demultiplexer	Average % ID (sd)	Average Alignment Length (bp) (sd)	Average Number of Clustered Reads (sd)	Average % Clustered Reads (sd)
100 reads per sample (100R)	MiniBar	99.99 (0.05)	421.05 (0.21)	97.5 (5.8)	97.50% (0.06)
	qcat	100 (0.00)	420.5 (0.86)	97.45 (6.01)	97.45% (0.06)
500 reads per sample (500R)	MiniBar	99.97 (0.11)	421.09 (0.43)	484.5 (35.77)	96.90% (0.07)
	qcat	100 (0.00)	420.82 (0.59)	483.68 (38.32)	96.73% (0.08)
5000 reads per sample (5KR)	MiniBar	99.88 (0.24)	421.18 (0.8)	4411.14 (916.69)	88.22% (0.18)
	qcat	99.95 (0.18)	420.41 (0.85)	4456.14 (939.87)	89.12% (0.19)

### 3.6. Validation of Sample Species Identity

The average sequence similarity between MinION consensus sequences and their corresponding Sanger sequence was highly accurate (>99.29% match) and remarkably consistent across sample type, extraction method, subset, and demultiplexer (Figure 4, Table 1). There was slightly more variation in sequence similarity across 5KR subsets, with the overall lowest percent sequence match (99.29%) obtained in these subsets for the cinnamon teal scat/Biomeme sample. This sample also had lower read cluster proportions (Figure 3) and the greatest loss in data after filtering (Figure S5).



**Figure 4.** The percent sequence similarity of MinION consensus to Sanger sequence from Blast for 100R, 500R, and 5KR subsets for each species. Samples are labeled by tissue type and extraction method (b = biomeme, c = chelex, q = qiagen). Points are linked by a grey line to show difference in values from demultiplexers. Overlapping areas in orange indicate similar results for MiniBar and qcat analyses. The horizontal dashed line is the 99% threshold for sequence similarity. Asterisks (\*) indicate samples with cinnamon teal read contamination.

The MinION consensus sequences from both MiniBar- and qcat-demultiplexed subsets extended into the Cytb primer region. We trimmed away the primers from both Sanger and MinION consensus

sequences for Mafft alignment of all samples. The cinnamon teal alignment had 99.8% pairwise identity and 97.2% identical sites ( $n = 84$  sequences), while the snow leopard alignment had 99.9% pairwise identity and 98.6% identical sites ( $n = 69$  sequences). The MinION consensus and Sanger sequences for each animal mainly differed at the ends of the sequences and at homopolymeric regions of varying lengths within the sequence (Table S7, Figure 5).



**Figure 5.** Screenshots of selected sections of the Mafft alignments for (A) snow leopard and (B) cinnamon teal showing nucleotide sites with differences between sequences in homopolymeric regions. Sanger sequences are listed above the black line and MinION consensus sequences below.

#### 4. Discussion

We demonstrate that a MinION-based DNA barcoding workflow can generate accurate consensus sequences from scat, hair, feather, and FFPE liver tissue samples, which are often considered challenging for molecular studies. The ability to use field-friendly DNA extraction protocols with these sample types will help to overcome logistical challenges, such as the need for cumbersome or expensive equipment, for molecular field research. The accuracy of our species identifications is on par with previous MinION DNA barcoding studies and pipelines [23–26,30]. For all tissue types, extraction methods, and subsets tested with our pipeline, we obtained high quality reads and a consensus sequence that matched >99.29% and at least 419/421 bp to the Sanger sequence for each sample. Although Oxford Nanopore’s goal is the “analysis of any living thing, by anyone, anywhere,” major barriers to its use are ease of sample processing, complicated data analysis, and cost. The results of our study help to reduce these barriers.

#### 4.1. Field-Friendly Protocols for Wildlife Samples Expands Conservation Applications with the MinION

We show that the Chelex and Biomeme extraction methods can be used to generate highly accurate MinION consensus sequences, similar to Qiagen extraction methods, even with low starting DNA concentrations. Our PCR amplicon purification and library prep protocols resulted in libraries of sufficient purity; cellular debris or contaminants present in the Chelex and Biomeme extracts did not affect sequencing of the Cytb amplicons. However, we have to caution that while PCR amplification resulted in high purity, direct sequencing of DNA extracts obtained using these methods is likely to be negatively affected by contaminants (e.g., by inhibiting library preparation or sequencing). ONT's protocols recommend using highly purified DNA for genomic sequencing. Although the field-friendly DNA extracts had low DNA concentrations overall, amplification was successful for all samples, including scat (known for containing PCR inhibitors), hair and feather (low DNA quantities), and FFPE tissue, from which DNA is generally difficult to amplify. Alternatively, ONT's Native Barcode Expansion kits can be used to ligate on indexes to the amplicon of interest.

Formalin can cause DNA fragmentation, cross-linking, subsequent sequence artifacts, and altered base pairs [16,17]. As artifacts are randomly distributed, they should not affect the final Sanger sequence if sufficient starting template is used [38,39]. Indeed, we accurately sequenced Qiagen-extracted DNA from FFPE samples, and further show that amplifiable DNA was successfully isolated from FFPE tissue using Chelex and Biomeme extraction methods.

#### 4.2. SAIGA: A DNA Barcoding Bioinformatics Pipeline for New MinION Users

We developed the SAIGA bioinformatics pipeline with a read clustering and consensus calling approach using software that were specifically designed for long-read and error-prone sequence data (isONclust, SPOA, Medaka). SAIGA performed successfully and consistently with as few as 100 reads per sample, allowing researchers to reduce sequencing time and cost per sample (e.g., multiplexing more samples). Like other studies investigating read coverage requirements, species identification accuracy still met our requirements but dropped slightly for the larger subset (5KR) [23,24]. Further, SAIGA options allow users to explore parameters and provide informative data quality checks and statistics throughout the pipeline. All software components are freely available, and the pipeline structure allows for integration of new software in the future.

Our results show that both qcat and MiniBar correctly demultiplex reads between samples in a sequence run and across multiple runs on a flow cell. Due to the very stringent demultiplexing parameters, the majority of raw data loss occurred during read assignment. More relaxed settings reduce raw read loss, but increase the chance of including incorrectly assigned reads or reads with higher sequencing error. Srivathsan et al. [26] and Maestri et al. [25] noted similar magnitudes of read loss with  $\approx 76\%$  and  $\approx 53.6\%$  of reads lost after demultiplexing, respectively; other MinION DNA barcoding publications have not reported this statistic. Despite the read loss, MiniBar- and qcat-demultiplexed reads performed well based on all our metrics for accurate species identification. Both demultiplexers tend to under-trim reads, which is preferred since potentially useful regions of the amplicon for distinguishing species are lost from over-trimmed reads. Although the consensus accuracy of qcat results was slightly higher than MiniBar results, we prefer MiniBar for its flexibility to analyze non-ONT index sequences. Customized indexes are less expensive than ONT indexes and can be lyophilized for field use.

Measuring the proportion of clustered filtered reads used for consensus sequence generation provides a benchmark for detecting sequencing error and potential contamination. For example, SAIGA created separate SPOA consensus sequence clusters for some samples even though these clusters produce the same species identification result. Lowering the sequence similarity threshold in cd-hit could force the sequences to form a single cluster. However, for the purpose of validating SAIGA, we used very stringent sequence similarity thresholds to reduce species identification bias from sequence error. Using this measure, we also show that SAIGA can handle low to medium amounts of laboratory contamination ( $\approx 4\%$ – $20\%$  reads of total subsample) from relatively distinct species in

samples without affecting final species identification since contaminant reads were successfully filtered out during the clustering process. Since contaminant teal reads had the correct indexes used for the three snow leopard samples, contamination likely occurred during library preparation rather than from mis-assignment of reads during demultiplexing. These snow leopard samples were either difficult to amplify during the Barcoding PCR (scat/Chelex) or had low recovery of indexed PCR product used in the sequencing run (hair/Biomeme and liver/Chelex). The contamination risk for these samples was likely exacerbated by the two-step PCR protocol and low starting DNA concentration and/or purity. Further development is needed to adapt this workflow and pipeline for mixed species samples, for which it may be more difficult to differentiate between true sample species and laboratory contaminants.

#### 4.3. Cost-Effective Strategies for Field Implementation

Each field-friendly method has its advantages and disadvantages. The Chelex method is cheap and the resin can be transported at room temperature, but requires heating equipment and the Chelex solution must be kept cool (4 °C) once prepared. The Biomeme kit is room temperature stable and self-contained. However, it is more expensive than both the Chelex resin and Qiagen kits (\$15/sample versus \$0.17 and \$3, respectively) and yielded lower DNA concentrations compared to the Qiagen kit.

We show that qcat and MiniBar can correctly assign reads to samples within and between runs, which reduces costs by allowing multiple sequence runs per flow cell. Future experiments can also scale up by sequencing more samples per flow cell because relatively few reads per sample are required for a consistent, accurate consensus (e.g., [26]). For the Cytb barcode amplified in this study, reads were sequenced at a rate of  $\approx 100,000$  reads per  $\approx 10$  min. Sufficient sequence data for species barcoding can therefore be obtained rapidly depending on the barcoding gene length and number of samples. We also reduced the volumes of the ONT PCR index per sample by 50% to lower costs and maximize the ONT kit.

## 5. Conclusions

Portable sequencing technology and field-friendly protocols have incredible potential to overcome institutional and geographical obstacles that impede genetic analyses in wildlife conservation and animal health. The methods described here provide an easy-to-follow workflow using field-friendly DNA extraction methods that can be used for preserved and non-invasively collected wildlife sample types to produce high-quality consensus sequences for species identification. Future studies are necessary to develop additional field-friendly protocols to further reduce the need for cold chain requirements, scale up sample processing, and tackle samples of mixed species, which will help to increase the opportunities for implementation.

**Supplementary Materials:** The following are available online at <http://www.mdpi.com/2073-4425/11/4/445/s1>, Figure S1: Gel electrophoresis comparison for tissue quality, Figure S2: Gel electrophoresis comparison for DNA extract quality; Figure S3: Flow cell active channels per sequence run, Figure S4: Basecall quality per sequence run, Figure S5: Proportion of demultiplexed reads after filtering for downstream clustering, Figure S6: Comparison of blast alignment length per analysis. Table S1: primer table, Table S2: DNA extraction concentrations, Table S3: MinION read statistics, Table S4: Experimental design information, Table S5: CD-HIT cluster results for proportions less than 60%, Table S6: Proportion of demultiplexed reads used for analysis, Table S7: MinION sequence error descriptions. Text S1: DNA extraction protocol details for each sample type and extraction method.

**Author Contributions:** Conceptualization, A.S., M.C.W.L., S.P., D.M. and T.A.S.; Data curation, A.S. and M.C.W.L.; Formal analysis, A.S. and M.C.W.L.; Funding acquisition, D.M. and T.A.S.; Investigation, A.S., M.C.W.L., S.P. and T.A.S.; Methodology, A.S., M.C.W.L., S.P. and T.A.S.; Project administration, T.A.S.; Resources, D.M. and T.A.S.; Software, M.C.W.L. and S.P.; Supervision, S.P., D.M. and T.A.S.; Validation, A.S. and M.C.W.L.; Visualization, M.C.W.L.; Writing—original draft, A.S. and M.C.W.L.; Writing—review and editing, A.S., M.C.W.L., S.P., D.M. and T.A.S. All authors have read and agreed to the published version of the manuscript.

**Funding:** Funding was provided by the G. Unger Vetlesen Foundation.

**Acknowledgments:** We thank Nina Vasiljevic and Rob Ogden for sharing their library preparation protocol and valuable discussions for our informatics pipeline, Batya Nightingale for lab assistance, and two anonymous reviewers for helpful comments.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

- Schlager, R.; Chiu, C.Y.; Miller, S.; Procop, G.W.; Weinstock, G. Validation of Metagenomic Next-Generation Sequencing Tests for Universal Pathogen Detection. *Arch. Pathol. Lab. Med.* **2017**, *141*, 776–786. [[CrossRef](#)] [[PubMed](#)]
- Gardy, J.L.; Loman, N.J. Towards a genomics-informed, real-time, global pathogen surveillance system. *Nat. Rev. Genet.* **2018**, *19*, 9–20. [[CrossRef](#)] [[PubMed](#)]
- Hobbs, C.A.D.; Potts, R.W.A.; Walsh, M.B.; Usher, J.; Griffiths, A.M. Using DNA Barcoding to Investigate Patterns of Species Utilisation in UK Shark Products Reveals Threatened Species on Sale. *Sci. Rep.* **2019**, *9*, 1–10. [[CrossRef](#)] [[PubMed](#)]
- Pardo, M.Á.; Jiménez, E.; Viðarsson, J.R.; Ólafsson, K.; Ólafsdóttir, G.; Danielsdóttir, A.K.; Pérez-Villareal, B. DNA barcoding revealing mislabeling of seafood in European mass caterings. *Food Control* **2018**, *92*, 7–16. [[CrossRef](#)]
- Galimberti, A.; Casiraghi, M.; Bruni, I.; Guzzetti, L.; Cortis, P.; Berterame, N.M.; Labra, M. From DNA barcoding to personalized nutrition: The evolution of food traceability. *Curr. Opin. Food Sci.* **2019**, *28*, 41–48. [[CrossRef](#)]
- Costa, F.O.; Carvalho, G.R. The Barcode of Life Initiative: Synopsis and prospective societal impacts of DNA barcoding of Fish. *Genom. Soc. Policy* **2007**, *3*, 29. [[CrossRef](#)]
- Ficetola, G.F.; Miaud, C.; Pompanon, F.; Taberlet, P. Species detection using environmental DNA from water samples. *Biol. Lett.* **2008**, *4*, 423–425. [[CrossRef](#)]
- Thomas, A.C.; Tank, S.; Nguyen, P.L.; Ponce, J.; Sinnesael, M.; Goldberg, C.S. A system for rapid eDNA detection of aquatic invasive species. *Environ. DNA* **2019**. [[CrossRef](#)]
- Marshall, H.D.; Ritland, K. Genetic diversity and differentiation of Kermode bear populations. *Mol. Ecol.* **2002**, *11*, 685–697. [[CrossRef](#)]
- Waits, L.P.; Paetkau, D. Noninvasive Genetic Sampling Tools for Wildlife Biologists: A Review of Applications and Recommendations for Accurate Data Collection. *J. Wildl. Manag.* **2005**, *69*, 1419–1433. [[CrossRef](#)]
- De Barba, M.; Miquel, C.; Boyer, F.; Mercier, C.; Rioux, D.; Coissac, E.; Taberlet, P. DNA metabarcoding multiplexing and validation of data accuracy for diet assessment: Application to omnivorous diet. *Mol. Ecol. Resour.* **2014**, *14*, 306–323. [[CrossRef](#)]
- Kohn, M.; Knauer, F.; Stoffella, A.; Schröder, W.; Pääbo, S. Conservation genetics of the European brown bear—A study using excremental PCR of nuclear and mitochondrial sequences. *Mol. Ecol.* **1995**, *4*, 95–104. [[CrossRef](#)]
- Rådström, P.; Knutsson, R.; Wolffs, P.; Lövenklev, M.; Löfström, C. Pre-PCR processing. *Mol. Biotechnol.* **2004**, *26*, 133–146. [[CrossRef](#)]
- Chaturvedi, U.; Tiwari, A.K.; Ratta, B.; Ravindra, P.V.; Rajawat, Y.S.; Palia, S.K.; Rai, A. Detection of canine adenoviral infections in urine and faeces by the polymerase chain reaction. *J. Virol. Methods* **2008**, *149*, 260–263. [[CrossRef](#)]
- Seimon, T.A.; Ayebare, S.; Sekisambu, R.; Muhindo, E.; Mitamba, G.; Greenbaum, E.; Menegon, M.; Pupin, F.; McAloose, D.; Ammazalorso, A.; et al. Assessing the Threat of Amphibian Chytrid Fungus in the Albertine Rift: Past, Present and Future. *PLoS ONE* **2015**, *10*, e0145841.
- Do, H.; Dobrovic, A. Sequence Artifacts in DNA from Formalin-Fixed Tissues: Causes and Strategies for Minimization. *Clin. Chem.* **2015**, *61*, 64–71. [[CrossRef](#)] [[PubMed](#)]
- Einaga, N.; Yoshida, A.; Noda, H.; Suemitsu, M.; Nakayama, Y.; Sakurada, A.; Kawaji, Y.; Yamaguchi, H.; Sasaki, Y.; Tokino, T.; et al. Assessment of the quality of DNA from various formalin-fixed paraffin-embedded (FFPE) tissues and the use of this DNA for next-generation sequencing (NGS) with no artifactual mutation. *PLoS ONE* **2017**, *12*, e0176280. [[CrossRef](#)]
- Hebert, P.D.N.; Ratnasingham, S.; de Waard, J.R. Barcoding animal life: Cytochrome c oxidase subunit 1 divergences among closely related species. *Proc. R. Soc. Lond. Ser. B Biol. Sci.* **2003**, *270*, S96–S99. [[CrossRef](#)]
- Valentini, A.; Pompanon, F.; Taberlet, P. DNA barcoding for ecologists. *Trends Ecol. Evol.* **2009**, *24*, 110–117. [[CrossRef](#)] [[PubMed](#)]
- Ip, C.L.C.; Loose, M.; Tyson, J.R.; de Cesare, M.; Brown, B.L.; Jain, M.; Leggett, R.M.; Eccles, D.A.; Zalunin, V.; Urban, J.M.; et al. MinION Analysis and Reference Consortium: Phase 1 data release and analysis. *F1000Reseach* **2015**, *4*, 1075. [[CrossRef](#)]
- Jain, M.; Tyson, J.R.; Loose, M.; Ip, C.L.C.; Eccles, D.A.; O’Grady, J.; Malla, S.; Leggett, R.M.; Wallerman, O.; Jansen, H.J.; et al. MinION Analysis and Reference Consortium: Phase 2 data release and analysis of R9.0 chemistry. *F1000Reseach* **2017**, *6*, 760. [[CrossRef](#)]

22. Krehenwinkel, H.; Pomerantz, A.; Prost, S. Genetic Biomonitoring and Biodiversity Assessment Using Portable Sequencing Technologies: Current Uses and Future Directions. *Genes* **2019**, *10*, 858. [[CrossRef](#)]
23. Pomerantz, A.; Peñafiel, N.; Arteaga, A.; Bustamante, L.; Pichardo, F.; Coloma, L.A.; Barrio-Amorós, C.L.; Salazar-Valenzuela, D.; Prost, S. Real-time DNA barcoding in a rainforest using nanopore sequencing: Opportunities for rapid biodiversity assessments and local capacity building. *GigaScience* **2018**, *7*, giy033. [[CrossRef](#)]
24. Krehenwinkel, H.; Pomerantz, A.; Henderson, J.B.; Kennedy, S.R.; Lim, J.Y.; Swamy, V.; Shoobridge, J.D.; Graham, N.; Patel, N.H.; Gillespie, R.G.; et al. Nanopore sequencing of long ribosomal DNA amplicons enables portable and simple biodiversity assessments with high phylogenetic resolution across broad taxonomic scale. *GigaScience* **2019**, *8*, giz006. [[CrossRef](#)]
25. Maestri, S.; Cosentino, E.; Paterno, M.; Freitag, H.; Garces, J.M.; Marcolungo, L.; Alfano, M.; Njunjić, I.; Schilthuizen, M.; Slik, F.; et al. A Rapid and Accurate MinION-Based Workflow for Tracking Species Biodiversity in the Field. *Genes* **2019**, *10*, 468. [[CrossRef](#)]
26. Srivathsan, A.; Hartop, E.; Puniamorthy, J.; Lee, W.T.; Kutty, S.N.; Kurina, O.; Meier, R. Rapid, large-scale species discovery in hyperdiverse taxa using 1D MinION sequencing. *BMC Biol.* **2019**, *17*, 96. [[CrossRef](#)]
27. Walsh, P.S.; Metzger, D.A.; Higuchi, R. Chelex 100 as a medium for simple extraction of DNA for PCR-based typing from forensic material. *Biotechniques* **1991**, *10*, 506–513. [[CrossRef](#)]
28. Singh, U.A.; Kumari, M.; Iyengar, S. Method for improving the quality of genomic DNA obtained from minute quantities of tissue and blood samples using Chelex 100 resin. *Biol. Proced. Online* **2018**, *20*, 12. [[CrossRef](#)]
29. Sepulveda, A.; Hutchins, P.; Massengill, R.; Dunker, K. Tradeoffs of a portable, field-based environmental DNA platform for detecting invasive northern pike (*Esox lucius*) in Alaska. *MBI* **2018**, *9*, 253–258. [[CrossRef](#)]
30. Srivathsan, A.; Baloglu, B.; Wang, W.; Tan, W.X.; Bertrand, D.; Ng, A.H.Q.; Boey, E.J.H.; Koh, J.J.Y.; Nagarajan, N.; Meier, R. A MinION™-based pipeline for fast and cost-effective DNA barcoding. *Mol. Ecol. Resour.* **2018**, *18*, 1035–1049. [[CrossRef](#)] [[PubMed](#)]
31. Verma, S.K.; Singh, L. Novel universal primers establish identity of an enormous number of animal species for forensic application. *Mol. Ecol. Notes* **2003**, *3*, 28–31. [[CrossRef](#)]
32. De Coster, W.; D’Hert, S.; Schultz, D.T.; Cruts, M.; Van Broeckhoven, C. NanoPack: Visualizing and processing long-read sequencing data. *Bioinformatics* **2018**, *34*, 2666–2669. [[CrossRef](#)]
33. Sahlin, K.; Medvedev, P. De novo clustering of long-read transcriptome data using a greedy, quality-value based algorithm. *bioRxiv* **2018**. [[CrossRef](#)]
34. Lee, C. Generating consensus sequences from partial order multiple sequence alignment graphs. *Bioinformatics* **2003**, *19*, 999–1008. [[CrossRef](#)]
35. Li, W.; Godzik, A. Cd-hit: A fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* **2006**, *22*, 1658–1659. [[CrossRef](#)]
36. Fu, L.; Niu, B.; Zhu, Z.; Wu, S.; Li, W. CD-HIT: Accelerated for clustering the next-generation sequencing data. *Bioinformatics* **2012**, *28*, 3150–3152. [[CrossRef](#)]
37. Altschul, S.F.; Gish, W.; Miller, W.; Myers, E.W.; Lipman, D.J. Basic local alignment search tool. *J. Mol. Biol.* **1990**, *215*, 403–410. [[CrossRef](#)]
38. Srinivasan, M.; Sedmak, D.; Jewell, S. Effect of Fixatives and Tissue Processing on the Content and Integrity of Nucleic Acids. *Am. J. Pathol.* **2002**, *161*, 1961–1971. [[CrossRef](#)]
39. Quach, N.; Goodman, M.F.; Shibata, D. In vitro mutation artifacts after formalin fixation and error prone translesion synthesis during PCR. *BMC Clin. Pathol.* **2004**, *4*, 1. [[CrossRef](#)]

