

## ORIGINAL RESEARCH

# NGSpeciesID: DNA barcode and amplicon consensus generation from long-read sequencing data

Kristoffer Sahlin<sup>1</sup>  | Marisa C. W. Lim<sup>2</sup>  | Stefan Prost<sup>3,4</sup> 

<sup>1</sup>Department of Mathematics, Science for Life Laboratory, Stockholm University, Stockholm, Sweden

<sup>2</sup>Department of Population Health and Reproduction, University of California, Davis, CA, USA

<sup>3</sup>LOEWE-Centre for Translational Biodiversity Genomics, Senckenberg, Frankfurt, Germany

<sup>4</sup>South African National Biodiversity Institute, National Zoological Garden, Pretoria, South Africa

**Correspondence**

Stefan Prost, LOEWE-Centre for Translational Biodiversity Genomics, Senckenberganlage 25, 60325 Frankfurt, Germany.  
Email: stefanprost.research@protonmail.com

**Abstract**

Third-generation sequencing technologies, such as Oxford Nanopore Technologies (ONT) and Pacific Biosciences (PacBio), have gained popularity over the last years. These platforms can generate millions of long-read sequences. This is not only advantageous for genome sequencing projects, but also advantageous for amplicon-based high-throughput sequencing experiments, such as DNA barcoding. However, the relatively high error rates associated with these technologies still pose challenges for generating high-quality consensus sequences. Here, we present NGSpeciesID, a program which can generate highly accurate consensus sequences from long-read amplicon sequencing technologies, including ONT and PacBio. The tool includes clustering of the reads to help filter out contaminants or reads with high error rates and employs polishing strategies specific to the appropriate sequencing platform. We show that NGSpeciesID produces consensus sequences with improved usability by minimizing preprocessing and software installation and scalability by enabling rapid processing of hundreds to thousands of samples, while maintaining similar consensus accuracy as current pipelines.

**KEYWORDS**

amplicon sequencing, DNA barcoding, sequence clustering, third-generation sequencing

## 1 | INTRODUCTION

We are in the middle of a biodiversity crisis, in which anthropogenic change is driving many species to extinction, often faster than they can be characterized (see, e.g., Ceballos et al., 2020). The identification of species in our environments is paramount to informing conservation policy and practice. The development of DNA barcoding (Hebert et al., 2003) was a major step toward large-scale characterizations of biodiversity. This technique utilizes amplification of standardized genetic regions to characterize species present within biological samples. Besides the documentation of biodiversity, this method and other amplicon sequencing technologies have been

widely used for monitoring of invasive species, detection of pathogens in environmental samples, and many other applications in taxonomy, medicine, or evolutionary biology (e.g., reviewed in Kress et al., 2015).

Third-generation sequencing is able to sequence millions of single molecules up to several Mbs in lengths (Jain et al., 2018). Currently, two platforms are readily available for DNA barcoding efforts, PacBio's Sequel II and ONT's MinION. These platforms offer the advantage of longer reads, at the cost of sequencing errors. While ONT's MinION still shows higher error rates >5% (Wick et al., 2018), the new PacBio HiFi mode allows for the generation of read with <1% error (Wenger et al., 2019), which will greatly improve the generation of

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2021 The Authors. *Ecology and Evolution* published by John Wiley & Sons Ltd

accurate DNA barcodes. Early on, researchers identified the potential of third-generation sequencing platforms for sequencing much longer DNA barcodes than previously possible (see, e.g., Tedersoo et al., 2018; Krehenwinkel, Pomerantz, Henderson, et al., 2019; Wurzbacher et al., 2019). Beside the longer amplicon length, ONT's MinION also offers the advantage that sequencing can be carried out almost anywhere in the world, due to its small size and affordability (reviewed in Krehenwinkel, Pomerantz, & Prost, 2019). While there has been a considerable software development effort to assemble high-quality amplicon consensus sequences from error-prone ONT MinION reads (see, e.g., Maestri et al., 2019; Seah et al., 2020; Srivathsan et al., 2019; reviewed in Krehenwinkel, Pomerantz, & Prost, 2019), only a few software solutions are available for PacBio-based DNA barcodes (see, e.g., Wurzbacher et al., 2019). To our knowledge, of these, only the pipeline presented in Wurzbacher et al. (2019) is able to handle both PacBio and ONT sequencing reads.

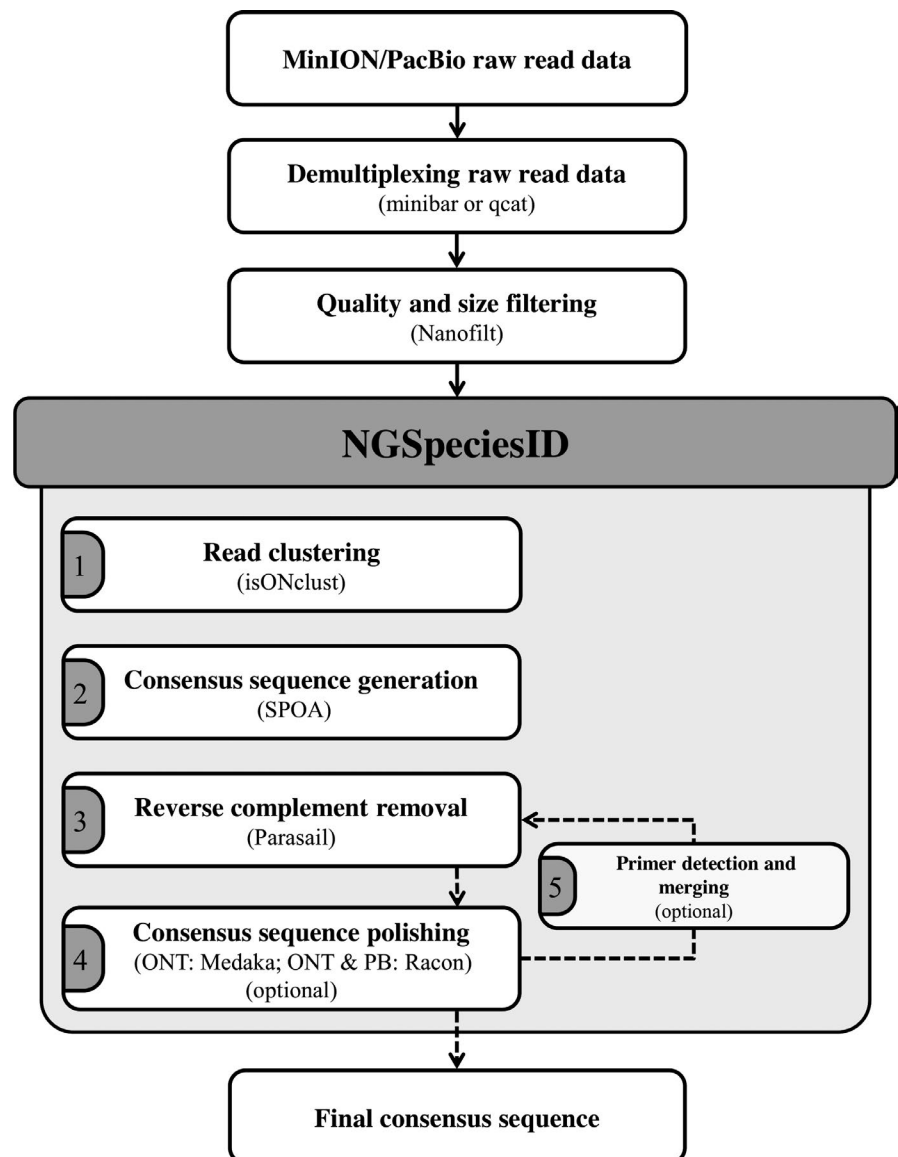
Here, we present NGSspeciesID a one-software solution for reconstructing high-quality amplicon consensus sequences for both

PacBio and ONT sequencing reads. We also investigate the performance of ONT's Medaka polishing software compared to Racon (Vaser et al., 2017) for MinION-based DNA barcoding. Compared to other programs, NGSspeciesID can be easily installed with conda, does not require any specific file name structures, can handle data from both third-generation sequencing types, includes different consensus polishing options, and only needs fastq files as input. We show that our tool produces consensus sequences of a similar quality than other software solutions, while reducing the burden to users by requiring little to no additional tools or data reformatting.

## 2 | SOFTWARE DESCRIPTION

NGSpeciesID is a program developed in python that wraps a set of tools for read clustering, consensus forming, and consensus polishing (Figure 1). It is a one-software solution and extension of the Saiga pipeline, we developed previously (Seah et al., 2020). It can be easily

**FIGURE 1** Steps involved in DNA barcode consensus calling of long-read data. The respective software tools used in the different steps are provided in brackets. In the first step, long-read data are usually demultiplexed. After demultiplexing, the reads are filtered for read length and quality. This step can also be carried out before demultiplexing if the respective amplicons do not differ in length. Next, consensus sequences for the individual read files can be generated using NGSspeciesID. If multiple read files need to be processed, NGSspeciesID can be run in a pipeline (see File S14). Within the tool, reads are first clustered according to similarity. Next, consensus sequences are generated for each cluster larger than an abundance threshold (default: >10% of all reads). In the third step, NGSspeciesID checks the generated consensus sequences for reverse complementarity. If consensus sequences are reverse complement, then the respective clusters are merged. In step four, the consensus sequences are polished using the reads from the respective clusters (this step is optional). In the last optional step, primers can be removed if this was not already carried out by the demultiplexing or basecalling tools. If primers are removed, NGSspeciesID will carry out steps 3–4 again



installed using the free and open-source Anaconda distribution. Briefly, NGSspeciesID clusters amplicon sequencing reads (in fastq format) and forms a consensus sequence for each cluster. Next, it merges reverse-complement clusters. Finally, the remaining consensus sequence(s) is/are polished. Optionally, the tool can also remove primer sequences from the consensus after the polishing step. In the following sections, we describe the workflow of NGSspeciesID, which is freely available at <https://github.com/ksahlin/NGSpeciesID>. For more details, see File S1 and Figure 1.

## 2.1 | Clustering of reads

NGSpeciesID first clusters input sequence reads based on expected sequence similarity for ONT or PacBio reads. Clustering is performed to remove any sequencing artifacts or contamination and assures that only similar reads from the relevant amplicon region will be considered when producing a polished amplicon. NGSspeciesID uses the isONclust clustering algorithm, which accounts for variable error rates within reads and is designed for both ONT or PacBio sequencing technologies. isONclust was recently shown to perform better than other clustering algorithms on both ONT or PacBio data (Sahlin & Medvedev, 2019).

## 2.2 | Forming draft consensus

Next, a draft consensus amplicon sequence is formed for each cluster that contains more reads than a specified proportion of the total reads (default: 10% of the total number of reads). The draft consensus sequences are formed with spoa (Vaser et al., 2017).

## 2.3 | Reverse-complement detection and removal

During the clustering step, reverse-complement reads from the amplicon region can produce two separate clusters of the same amplicon. In order to assure only one sequence per amplicon, NGSspeciesID detects and merges any consensus sequences classified as reverse-complement sequences with respect to each other using pairwise alignment with parasail (Daily, 2016). To do so, all consensus sequences are aligned to each other and any two sequences with alignment identity over a parameter (default: 90%) are merged. The original reads that were used to generate the two consensus sequences are combined to increase coverage. Finally, all draft consensus sequences passing this step, with the original reads, are sent to the polishing step.

## 2.4 | Polishing

The remaining consensus sequences are polished with either Medaka (<https://github.com/nanoporetech/medaka>) or Racon

(Vaser et al., 2017). In this step, the original reads are first mapped back to the consensus sequence (reference sequence). The reference is then corrected using sequence information from the multiple reads mapped. The polished consensus sequences are the final output of NGSspeciesID.

## 2.5 | Primer detection and removal

Many basecalling and demultiplexing tools do not remove primers from the amplicon sequences (but see Minibar (Krehenwinkel, Pomerantz, Henderson, et al., 2019)). NGSspeciesID, therefore, implements an optional primer removal step by searching the forward and reverse complement of each primer within a window at each end of the read. This step is carried out for the polished sequences to improve the detection of the priming sites. If no primer is found, the polished consensus sequence(s) remain the final output of NGSspeciesID. If primer(s) have been detected and trimmed, NGSspeciesID reruns the reverse-complement removal and polishing steps to identify any remaining redundant consensus sequences that were not removed due to primers.

## 3 | USE CASES AND COMPARISON TO OTHER TOOLS

We tested our software on publicly available data from Maestri et al. (2019) and Wurzbacher et al. (2019), and compared the accuracy of respective consensus sequences generated in the two studies to those reconstructed with NGSspeciesID. To measure accuracy, we aligned the consensus sequences to the respective Sanger sequence using BLAST (Altschul et al., 1990) and calculated accuracy as the sum of all matches in the alignment divided by the alignment length. We chose the two software solution (Mothur (Schloss et al., 2009) and Consension (Wurzbacher et al., 2019)) presented in Wurzbacher et al. (2019) for our comparison as it is currently, to our knowledge, the only one that can be used with both PacBio and ONT sequencing reads. We further compared our result to the ONTrack software (Maestri et al., 2019) developed for ONT data specifically. In both comparisons, we carried out polishing within NGSspeciesID using Medaka (<https://github.com/nanoporetech/Medaka>) and Racon (Vaser et al., 2017).

### 3.1 | Comparison to Mothur + Consension

We randomly selected five out of the 61 fungi datasets from Wurzbacher et al. (2019), ranging from 201 to 447 reads per dataset (Table S1). These cover five fungi species of the genus *Inocybe* for ribosomal DNA (rDNA) and the full ribosomal tandem repeat region (TR). We provide alignments of the corresponding Sanger sequences with our consensus sequences in the Files S2–S6. In their approach, Wurzbacher et al. (2019) first perform operational

taxonomic unit (OTU) clustering on the read data using Mothur (Schloss et al., 2009). Next, they create consensus sequences using Consension (Wurzbacher et al., 2019).

In general, we see that for both ONT and PacBio data NGSspeciesID and the Mothur + Consension software perform equally well, generating consensus sequences with 98.6% to 100% accuracy (Table 1). In three out of the five cases, the two pipelines produced consensus sequences with the same accuracy, while in one case each software slightly outperformed the other (Table 1). Medaka polishing outperformed Racon polishing in four out of five cases (Table 1).

### 3.2 | Comparison to ONTrack

Next, we compared the performance of NGSspeciesID to the pipeline ONTrack from Maestri et al. (2019). ONTrack first clusters all reads using VSEARCH (Rognes et al., 2016), then randomly selects 200 reads, aligns those with Mafft (Katoh & Standley, 2013), calls the consensus with EMBOSS cons (<http://emboss.sourceforge.net/apps/cvs/emboss/apps/cons.html>), and lastly carries out polishing with 200 randomly selected reads using Nanopolish (<https://github.com/jts/nanopolish>). We generated consensus sequences for all seven DNA barcodes from Maestri et al. (2019), which comprise *cytochrome C oxidase subunit 1* (COI) sequences of two snails and five beetles (Table S1). We provide the respective alignments in the Files S7–S13).

Previously, Krehenwinkel, Pomerantz, Henderson, et al. (2019) showed that consensus accuracy can decrease when too many reads (in the realm of a few hundred reads, depending on the error rate of the individual reads) are selected for the consensus generation, likely due to an increase in the signal to noise ratio. We thus randomly subsampled 300 reads using seqtk (command: `seqtk sample -s 1234 reads.fastq 300 > reads_subsample.fastq`; <https://github.com/lh3/seqtk>), a number which has been shown to work well with Nanopore data (Krehenwinkel, Pomerantz, Henderson, et al., 2019). We see that the consensus quality is comparable between the two tools (Table 2), with accuracy of 99.8% to 100%. In five out of the seven DNA barcode sets, both tools performed equally well, while in

one each the two tools outperformed each other, however, differing by only 1 base pair (Table 2).

### 3.3 | Mixed samples

We tested NGSspeciesID's performance on mixed samples in silico by combining 300 reads of each of the seven barcodes from Maestri et al. (2019). To do so, we set the cluster abundance ratio to 5% (--abundance\_ratio 0.05). We recovered seven consensus sequences corresponding to the seven DNA barcodes, ranging from 99.3% to 100% similarity to the corresponding Sanger sequence (Table 2). In four out of the seven cases, we recovered the same percentage similarity to the Sanger sequence in the mixed analysis as in the respective single barcode processing. In three cases, the accuracy was slightly lower with two and four basepair differences, respectively.

## 4 | DISCUSSION

### 4.1 | Consensus quality

Here, we present NGSspeciesID, an easy-to-use, one-software solution for the generation of high-quality consensus sequences for the long-read sequencing technologies from ONT and PacBio. We compared NGSspeciesID against results obtained with Mothur + Consension and ONTrack. In general, all three software solutions produced consensus sequences of a very high quality, reaching 99%–100% accuracy in almost all cases. We show that NGSspeciesID performs comparably to the other tools. Throughout all comparisons, we see that consensus sequences based on ONT data polished with Racon usually show lower percent similarities to the Sanger sequence than consensus sequences polished with Medaka. NGSspeciesID carries out 2 rounds of Racon polishing by default. Increasing or decreasing the number of rounds might increase the consensus quality. We chose Medaka as the default error corrector in NGSspeciesID as it includes up to date error models. We did not include an option to use Nanopolish in NGSspeciesID, which is used in ONTrack, as this tool requires fast5 files, which are often

**TABLE 1** Percent similarity to the respective Sanger sequence for the datasets 17075, 17078, 16416, 16427, and 16483 from Wurzbacher et al. (2019)

SampleID	17075	17078	16416	16427	16483
NGSpeciesID					
ONT Medaka	<b>98.6% (10/726)</b>	<b>99.2% (6/741)</b>	99.5% (4/731)	<b>99.6% (3/790)</b>	<b>100% (0/709)</b>
ONT Racon	98.5% (11/726)	99.1% (7/741)	99.7% (2/731)	99.1% (7/790)	99.9% (1/709)
PB Racon	<b>98.6% (10/726)</b>	99.1% (7/741)	<b>99.9% (1/731)</b>	<b>99.6% (3/790)</b>	<b>100% (0/709)</b>
Mothur + Consension					
ONT	<b>98.6% (10/726)</b>	<b>99.2% (6/741)</b>	<b>99.9% (1/731)</b>	99.5% (4/790)	<b>100% (0/709)</b>
PB	<b>98.6% (10/726)</b>	<b>99.2% (6/741)</b>	99.7% (2/731)	<b>99.6% (3/790)</b>	<b>100% (0/709)</b>

Note: The highest similarity scores are highlighted in bold. The numbers in the brackets provide the amount of mismatches to the Sanger reference and the length of the reference sequence.

**TABLE 2** Percent similarity to the respective Sanger sequence for the datasets B1 to BC7 from Maestri et al. (2019)

SampleID	BC1 <sup>a</sup>	BC2	BC3	BC4 <sup>a</sup>	BC5	BC6 <sup>a</sup>	BC7 <sup>a</sup>
NGSpeciesID							
ONT Medaka	<b>100% (0/651)</b>	<b>100% (0/658)</b>	99.9% (1/649)	<b>100% (0/606)</b>	<b>100% (0/658)</b>	<b>99.8% (1/576)</b>	<b>100% (0/536)</b>
ONT Racon	99.5% (3/651)	99.5% (3/658)	98.9% (7/649)	99.2% (5/606)	99.8% (1/658)	99.7% (2/576)	99.4% (3/536)
ONTrack							
ONT	99.9% (1/651)	<b>100% (1<sup>b</sup>/658)</b>	<b>100% (2<sup>b</sup>/649)</b>	<b>100% (0/606)</b>	<b>100% (2<sup>b</sup>/658)</b>	<b>99.8% (1/576)</b>	<b>100% (0/536)</b>
Mixed							
NGSpeciesID							
ONT Medaka	<b>100% (0/651)</b>	<b>100% (0/658)</b>	99.7% (2/649)	99.3% (4/606)	<b>100% (0/658)</b>	<b>99.8% (1/576)</b>	99.6% (2/536)

Note: For the mixed samples, 300 reads of each of the seven DNA barcodes were combined into a single file, from which NGSspeciesID generated multiple consensus sequences. NGSspeciesID was run using Medaka polishing.

<sup>a</sup>Here, the Sanger sequence from Maestri et al. (2019) was shorter than the expected fragment length and all the consensus sequences. In these cases, we only calculated the percentage similarity for the region covered by the respective Sanger sequence.

<sup>b</sup>The consensus sequences from Maestri et al. (2019) are missing one or two bases at the start, which could be due to a consensus calling error, or deletion of one additional base during the primer removal. For the percentage accuracy, we assumed them to be incorrectly trimmed. The highest similarity scores are highlighted in bold. The numbers in the brackets provide the amount of mismatches to the Sanger reference and the length of the reference sequence.

not available for published Oxford Nanopore data. Furthermore, it requires preprocessing to generate the appropriate header structure in the corresponding fastq files, which makes it much more time consuming to use.

As the generation of consensus sequences for DNA barcoding takes only a few seconds for each sample (depending on the number of reads), we did not compare run times between the different pipelines.

## 4.2 | Easy use

NGSpeciesID was designed to be straightforward to use. It works on individual read files, outputted either directly from the basecalling or after demultiplexing (e.g., using Minibar (Krehenwinkel, Pomerantz, Henderson, et al., 2019) or qcat (<https://github.com/nanoporetech/qcat>)), but can quickly be adjusted to run in a loop over multiple fastq files using a bash script (see File S14). It only requires fastq files as input. In contrast, ONTrack requires the input reads in three formats (fast5, fasta and fastq), which requires additional preprocessing of the sequencing data. Furthermore, NGSspeciesID allows fastq files to have any naming structure, thus making it easy for the user to run and to identify samples and replicates. This saves time on preprocessing of the read data compared to other software solutions.

NGSpeciesID employs quality filtering of the reads based on read phred scores. However, we recommend also removing reads much shorter or longer than the intended target, which often represent chimeras or contaminations using NanoFilt (De Coster et al., 2018) before running NGSspeciesID. While our tool can handle unfiltered data, this might result in the generation of multiple consensus sequences. NGSspeciesID also offers the option to remove priming sites from the amplicon sequences. As many universal primers include ambiguity codes, primer regions can potentially include incorrect

bases and should thus be removed. We further found that primer regions can cause issues for the reverse-complement matching. We thus included an additional reverse-complement matching step after primer removal, in case NGSspeciesID outputs multiple consensus sequences. Our tool outputs multiple consensus sequences in case the clustering results in multiple clusters over a certain percentage of the total reads (by default this is set to 10%). Each consensus sequence is only polished with the corresponding reads from the clustering. This feature is very useful as it allows the user to explore potential contaminant reads or mixed samples through the generating of multiple consensus sequences.

NGSpeciesID and the Mothur + Consension software solution both can handle ONT and PacBio long-read data. While both tools produce consensus sequences of similar accuracy, Mothur + Consension requires an in-depth knowledge of the pipeline requiring (a) preprocessing of the input data, (b) individual components of the pipeline to be run separately, and (c) has parameter settings that are difficult to interpret, while NGSspecies is designed to be user friendly and packaged as a one command solution.

## 4.3 | Mixed samples

While NGSspeciesID was not designed specifically for metabarcoding data, the flexibility of the algorithmic steps in the pipeline enables the tool to handle mixed samples if they are sufficiently divergent. We recovered seven consensus sequences corresponding to the seven DNA barcodes pooled in the mixed sample analysis. NGSspeciesID generated highly accurate consensus sequences for all barcodes, ranging from 99.2% to 100%. For the mixed sample test, we adjusted the read abundance ratio for the clusters to 5%, since the seven barcodes at equal abundance are each present in only 14% of the reads in the sample. Therefore, the default abundance cutoff of 10% would require 210



out of the 300 reads to be used per cluster, which might not be the case. Three out of seven barcodes showed a slightly lower consensus accuracy than in the respective single species analysis, which is likely due to the presence of some reads from other barcodes in the clusters that might have affected the polishing accuracy, and the random selection of the 300 reads for each barcode (as individual read error rates can differ). We expect some cross-contamination (reads assigned to the wrong cluster), especially for closely related species. However, this should improve with the continued improvement of third-generation sequencing read accuracy. This experiment shows that NGSpeciesID, even though it was not developed for mixed samples, can recover highly accurate consensus sequences from metabarcoding data if the samples are sufficiently divergent. However, its performance on metabarcoding data will need to be investigated separately with mock datasets of varying ratios and sample relationships (to see which taxonomic divergences are needed for effective separation of reads from related species).

## 5 | CONCLUSION AND FUTURE DIRECTIONS

We present NGSpeciesID, an easy-to-use and flexible one-software solution to generate high-quality consensus sequences for both ONT and PacBio sequencing data. It performs equally well as other pipelines and software solutions tested here, but offers advanced usability as it is simple to use and does not require preprocessing of the data before running. Portable devices such as the inexpensive MinION sequencer have started to democratize the process of molecular biodiversity monitoring (see, e.g., Krehenwinkel, Pomerantz, and Prost (2019)). Here, we add to this, by the development of a simple to install and run bioinformatic software that should further enable students and citizen scientists without a formalized bioinformatic training to carry out biodiversity monitoring and assessment studies.

### ACKNOWLEDGMENTS

We thank Tilman Schell and Aaron Pomerantz for their valuable comments on the paper draft, and Christian Wurzbacher and Simone Maestri for providing the consensus and Sanger sequences from their study. The authors declare no conflict of interest.

### CONFLICT OF INTEREST

The authors have declared that no competing interests exist.

### AUTHOR CONTRIBUTION

**Kristoffer Sahlin:** Conceptualization (Equal), Software (Equal), Writing-original draft (Equal), Writing-review & editing (Equal)  
**Marisa C. W. Lim:** Conceptualization (equal); Software (equal); Writing-original draft (equal); Writing-review & editing (equal).  
**Stefan Prost:** Conceptualization (Lead), Methodology (Equal), Project administration (Lead), Software (Equal), Supervision (Lead), Writing-original draft (Lead), Writing-review & editing (Lead)

### ETHICS STATEMENT

The presented study only used publicly available data.

### DATA AVAILABILITY STATEMENT

We did not generate sequencing read data within this study. GenBank accession numbers for all samples used in this study along with the citations of the papers they were published in are provided in Table S1. The software along with example read data can be found on <https://github.com/ksahlin/NGSpeciesID>.

### ORCID

Kristoffer Sahlin  <https://orcid.org/0000-0001-7378-2320>

Marisa C. W. Lim  <https://orcid.org/0000-0003-2097-8818>

Stefan Prost  <https://orcid.org/0000-0002-6229-3596>

### REFERENCES

- Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local alignment search tool. *Journal of Molecular Biology*, 215, 403–410. [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2)
- Ceballos, G., Ehrlich, P. R., & Raven, P. H. (2020). Vertebrates on the brink as indicators of biological annihilation and the sixth mass extinction. *Proceedings of the National Academy of Sciences of the United States of America*, 117, 13596–13602. <https://doi.org/10.1073/pnas.1922686117>
- Daily, J. (2016). Parasail: SIMD C library for global, semi-global, and local pairwise sequence alignments. *BMC Bioinformatics*, 17, 81. <https://doi.org/10.1186/s12859-016-0930-z>
- De Coster, W., D'Hert, S., Schultz, D. T., Cruts, M., & Van Broeckhoven, C. (2018). NanoPack: Visualizing and processing long-read sequencing data. *Bioinformatics*, 34, 2666–2669. <https://doi.org/10.1093/bioinformatics/bty149>
- Hebert, P. D. N., Ratnasingham, S., & de Waard, J. R. (2003). Barcoding animal life: Cytochrome c oxidase subunit 1 divergences among closely related species. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 270, S96–S99. <https://doi.org/10.1098/rsbl.2003.0025>
- Jain, M., Koren, S., Miga, K. H., Quick, J., Rand, A. C., Sasani, T. A., Tyson, J. R., Beggs, A. D., Dilthey, A. T., Fiddes, I. T., Malla, S., Marriott, H., Nieto, T., O'Grady, J., Olsen, H. E., Pedersen, B. S., Rhie, A., Richardson, H., Quinlan, A. R., ... Loose, M. (2018). Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nature Biotechnology*, 36, 338–345. <https://doi.org/10.1038/nbt.4060>
- Katoh, K., & Standley, D. M. (2013). MAFFT Multiple Sequence Alignment Software Version 7: Improvements in performance and usability. *Molecular Biology and Evolution*, 30, 772–780. <https://doi.org/10.1093/molbev/mst010>
- 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., & Prost, S. (2019). Nanopore sequencing of long ribosomal DNA amplicons enables portable and simple biodiversity assessments with high phylogenetic resolution across broad taxonomic scale. *Gigascience*, 8, giz006. <https://doi.org/10.1093/gigascience/giz006>
- Krehenwinkel, H., Pomerantz, A., & Prost, S. (2019). Genetic Biomonitoring and Biodiversity Assessment Using Portable Sequencing Technologies: Current uses and future directions. *Genes*, 10, 858. <https://doi.org/10.3390/genes10110858>
- Kress, W. J., García-Robledo, C., Uriarte, M., & Erickson, D. L. (2015). DNA barcodes for ecology, evolution, and conservation. *Trends in Ecology & Evolution*, 30, 25–35. <https://doi.org/10.1016/j.tree.2014.10.008>
- Maestri, S., Cosentino, E., Paterno, M., Freitag, H., Garces, J. M., Marcolungo, L., Alfano, M., Njunjić, I., Schilthuizen, M., Slik, F.,

- Menegon, M., Rossato, M., & Delledonne, M. (2019). A Rapid and Accurate MinION-Based workflow for tracking species biodiversity in the field. *Genes*, *10*, 468. <https://doi.org/10.3390/genes10060468>
- Rognes, T., Flouri, T., Nichols, B., Quince, C., & Mahé, F. (2016). VSEARCH: A versatile open source tool for metagenomics. *PeerJ*, *4*, e2584. <https://doi.org/10.7717/peerj.2584>
- Sahlin, K., & Medvedev, P. (2019). De Novo clustering of long-read transcriptome data using a greedy, quality-value based algorithm. In L. J. Cowen (Ed.), *Research in Computational Molecular Biology, Lecture Notes in Computer Science* (pp. 227–242). Springer International Publishing.
- Schloss, P. D., Westcott, S. L., Ryabin, T., Hall, J. R., Hartmann, M., Hollister, E. B., Lesniewski, R. A., Oakley, B. B., Parks, D. H., Robinson, C. J., Sahl, J. W., Stres, B., Thallinger, G. G., Horn, D. J. V., & Weber, C. F. (2009). Introducing mothur: Open-Source, Platform-Independent, Community-supported software for describing and comparing microbial communities. *Applied and Environment Microbiology*, *75*, 7537–7541. <https://doi.org/10.1128/AEM.01541-09>
- Seah, A., Lim, M. C. W., McAloose, D., Prost, S., & Seimon, T. A. (2020). MinION-Based DNA barcoding of preserved and non-invasively collected wildlife samples. *Genes*, *11*, 445. <https://doi.org/10.3390/genes11040445>
- Srivathsan, A., Hartop, E., Puniamoorthy, J., Lee, W. T., Kutty, S. N., Kurina, O., & Meier, R. (2019). Rapid, large-scale species discovery in hyperdiverse taxa using 1D MinION sequencing. *BMC Biology*, *17*, 96. <https://doi.org/10.1186/s12915-019-0706-9>
- Tedersoo, L., Tooming-Klunderud, A., & Anslan, S. (2018). PacBio metabarcoding of fungi and other eukaryotes: Errors, biases and perspectives. *New Phytologist*, *217*, 1370–1385. <https://doi.org/10.1111/nph.14776>
- Vaser, R., Sović, I., Nagarajan, N., & Šikić, M. (2017). Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Research*, *27*, 737–746. <https://doi.org/10.1101/gr.214270.116>
- Wenger, A. M., Peluso, P., Rowell, W. J., Chang, P.-C., Hall, R. J., Concepcion, G. T., Ebler, J., Functammasan, A., Kolesnikov, A., Olson, N. D., Töpfer, A., Alonge, M., Mahmoud, M., Qian, Y., Chin, C.-S., Phillippy, A. M., Schatz, M. C., Myers, G., DePristo, M. A., ... Hunkapiller, M. W. (2019). Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nature Biotechnology*, *37*, 1155–1162. <https://doi.org/10.1038/s41587-019-0217-9>
- Wick, R. R., Judd, L. M., & Holt, K. E. (2018). Deepbiner: Demultiplexing barcoded Oxford Nanopore reads with deep convolutional neural networks. *PLoS Computational Biology*, *14*, <https://doi.org/10.1371/journal.pcbi.1006583>
- Wurzbacher, C., Larsson, E., Bengtsson-Palme, J., den Wyngaert, S. V., Svantesson, S., Kristiansson, E., Kagami, M., & Nilsson, R. H. (2019). Introducing ribosomal tandem repeat barcoding for fungi. *Molecular Ecology Resources*, *19*, 118–127. <https://doi.org/10.1111/1755-0998.12944>

## SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section.

**How to cite this article:** Sahlin K, Lim MCW, Prost S. NGSspeciesID: DNA barcode and amplicon consensus generation from long-read sequencing data. *Ecol Evol*. 2021;11:1392–1398. <https://doi.org/10.1002/ece3.7146>