

Minimizing polymerase biases in metabarcoding

Ruth V. Nichols¹ | Christopher Vollmers² | Lee A. Newsom³ | Yue Wang⁴ |
Peter D. Heintzman^{1,5} | McKenna Leighton² | Richard E. Green² | Beth Shapiro¹

¹Department of Ecology and Evolutionary Biology, University of California Santa Cruz, Santa Cruz, California

²Department of Biomolecular Engineering, University of California Santa Cruz, Santa Cruz, California

³Department of Social Sciences, Flagler College, St. Augustine, Florida

⁴Department of Geography, University of Wisconsin-Madison, Madison, Wisconsin

⁵Tromsø University Museum, UiT – The Arctic University of Norway, Tromsø, Norway

Correspondence

Ruth V. Nichols and Beth Shapiro,
Department of Ecology and Evolutionary Biology, University of California Santa Cruz, Santa Cruz, CA.

E-mails: ruthnichols@gmail.com;
beth.shapiro@gmail.com

Funding information

University of California Office of the President, Grant/Award Number: 20160713SC; National Science Foundation, Grant/Award Number: ARC-1203990; Gordon and Betty Moore Foundation, Grant/Award Number: GBMF-3804

Abstract

DNA metabarcoding is an increasingly popular method to characterize and quantify biodiversity in environmental samples. Metabarcoding approaches simultaneously amplify a short, variable genomic region, or “barcode,” from a broad taxonomic group via the polymerase chain reaction (PCR), using universal primers that anneal to flanking conserved regions. Results of these experiments are reported as *occurrence* data, which provide a list of taxa amplified from the sample, or *relative abundance* data, which measure the relative contribution of each taxon to the overall composition of amplified product. The accuracy of both occurrence and relative abundance estimates can be affected by a variety of biological and technical biases. For example, taxa with larger biomass may be better represented in environmental samples than those with smaller biomass. Here, we explore how polymerase choice, a potential source of technical bias, might influence results in metabarcoding experiments. We compared potential biases of six commercially available polymerases using a combination of mixtures of amplifiable synthetic sequences and real sedimentary DNA extracts. We find that polymerase choice can affect both occurrence and relative abundance estimates and that the main source of this bias appears to be polymerase preference for sequences with specific GC contents. We further recommend an experimental approach for metabarcoding based on results of our synthetic experiments.

KEYWORDS

bias, eDNA, environmental DNA, metabarcoding, soil, trnL P6 loop

1 | INTRODUCTION

Metabarcoding, which is erroneously described as barcoding or metagenomics in some literature, is the technique in which a universal primer pair is used to amplify multiple templates from a mixture of many different taxa or haplotypes. Metabarcoding is often used in conjunction with environmental DNA (eDNA), or DNA that is collected from environmental sources such as water, sediment, air and faeces (Deiner et al., 2017). Metabarcoding is an increasingly popular tool in ecological and palaeoecological research, mainly due to its simplicity and low cost. eDNA can be used, for example, to characterize biodiversity of a particular taxonomic group (Ushio et al., 2017) or to estimate the ranges of rare, extinct or cryptic species

(Haile et al., 2009; Jerde, Mahon, Chadderton, & Lodge, 2011; Pedersen et al., 2016; Rees, Baker, Gardner, Maddison, & Gough, 2017). Additionally, metabarcoding has been used to calculate differences in haplotype or allele frequency between populations of the same species (Sigsgaard et al., 2016) and to link changes in community composition over time to climatic shifts (Haile et al., 2007; Willerslev et al., 2003, 2007, 2014). These latter examples analyse both the *occurrence* and *relative abundance* of each unique sequence in the amplification product, where abundance is estimated as the proportion of the total number of sequences generated matching each taxon or haplotype.

While metabarcoding is a promising approach to characterize biodiversity both quickly and inexpensively, few studies have

validated the method experimentally by, for example, testing the extent to which the true community or population is reconstructed. It is generally accepted that taxon occurrence can be inferred via metabarcoding, provided that a sufficient number of PCR replicates—amplifying DNA multiple times from the same soil extract using the same amplification conditions—are performed (Piñol, Mir, Gomez-Polo, & Agustí, 2015; Shaw et al., 2016) and false positives have been accounted for (Lahoz-Monfort, Guillera-Aroita, & Tingley, 2016). The first eDNA metabarcoding studies used replication (Cooper & Poinar, 2000), where DNA extraction and amplification were both replicated, to help confirm their results (Willerslev et al., 2003), but many subsequent studies did not replicate experiments (Soininen et al., 2009; Sønstebo et al., 2010; Valentini et al., 2009). After a detailed exploration of the utility of replication in metabarcoding (Darling & Mahon, 2011), the use of replication increased, but the number of replicates performed per experiment varied widely. Most studies used between two and five PCR replicates per sample (Andersen et al., 2012; De Barba et al., 2014; Jørgensen et al., 2012; Willerslev et al., 2014) and some as many as eight (Giguët-Covex et al., 2014). Recently, the use of site occupancy models has been proposed as a tool to estimate how many replicates are needed; with most recommendations ranging from six to 12 replicates per sample (Ficetola et al., 2015; Lahoz-Monfort et al., 2016; Schmidt, Kéry, Ursenbacher, Hyman, & Collins, 2013), depending on the number and abundance of rare taxa. Another approach to estimate the amount of replication required is rarefaction, whereby the number of new taxa identified per replicate PCR is used to estimate the probability that most rare taxa have been recovered (Hsieh, Ma, & Chao, 2016; Sanders, 1968).

Whether relative abundance can be estimated accurately from metabarcoding data is a more contentious issue. Some researchers routinely interpret the relative abundance of sequences post-PCR as indicative of real relative biomass estimates (Kowalczyk, Taberlet, Kaminski, & Wojcik, 2011; Niemyer, Epp, Stoof-Leichsenring, Pestryakova, & Herzsuh, 2017; Willerslev et al., 2014). Others argue against this approach, citing challenges that include differential DNA degradation, different primer binding efficiencies and sequencing errors as confounding factors that might influence the utility of relative abundance data collected from metabarcoding loci (Deagle, Thomas, Shaffer, Trites, & Jarman, 2013; Deagle et al., 2007; Marcelino & Verbruggen, 2016; Pawluczyk et al., 2015; Piñol et al., 2015).

Biases that might influence the likelihood of a taxon being detected during metabarcoding can be both biological and technical in origin. Biological differences include organism size, seasonal presence and senescence, preservation and dispersal strategy, among others. Larger taxa, taxa that are present year-round or taxa whose DNA is readily transported across long distances by wind or water, may be more likely to be observed in environmental samples than smaller, seasonal and sedentary taxa (Andersen et al., 2012; Barnes & Turner, 2016; Buxton, Groombridge, Zakaria, & Griffiths, 2017; Dunn, Priestley, Herraiz, Arnold, & Savolainen, 2017; Hemery, Politano, & Henkel, 2017; Rees et al., 2017). Even when the same number of cells is present in an environmental sample, the starting copy

number of target loci may vary between taxa and tissue type. Chloroplast DNA, for example, is a common target for metabarcoding, but can differ in copy number between taxa, individuals and cell tissue types within the same plant (Morley & Nielsen, 2016). Taphonomic factors may also influence DNA preservation, for example by affecting the rate of degradation. Lignified structures in plants may slow the rate of DNA degradation (Yoccoz et al., 2012), as may anoxic environments (Corinaldesi, Barucca, Luna, & Dell'Anno, 2011). In some environments, soil leaching and postdepositional mixing may move DNA up or down sediment columns or horizontally over space (Andersen et al., 2012; Anderson-Carpenter et al., 2011; Pedersen et al., 2015; Rawlence et al., 2014).

Technical biases can be introduced during DNA extraction and PCR amplification. DNA extraction protocols can be more or less optimized for soil chemistry, which can influence the extent to which DNA is recovered (Zielińska et al., 2017). Soils rich in clays or humic acids may bind DNA, for example, reducing DNA recovery (Direito, Marees, & Röling, 2012). PCR is a highly stochastic process, which is further complicated by the presence of variable templates, with many opportunities for the introduction of bias (Aird et al., 2011; Pinto & Raskin, 2012; Polz & Cavanaugh, 1998; Suzuki & Giovannoni, 1996). Although the universal primers used in metabarcoding are designed to anneal to conserved genomic regions, slight variation in binding site sequences may affect primer binding efficiency, resulting in bias (Elbrecht & Leese, 2015; Piñol et al., 2015). For example, Fahner, Shokralla, Baird, and Hajibabaei (2016) used four plant-specific primers to infer community composition from the same soil samples and found that each primer pair produced a different result. This result may also be related to amplicon length whereby shorter amplicons amplify more readily than longer amplicons. Template secondary structures can also bias PCR when molecules with secondary structures bind to themselves and inhibit their own amplification. In addition, templates with suboptimal GC contents can be disfavoured during amplification, although some polymerases are known to have reduced GC bias and additives such as dimethyl sulphoxide (DMSO) for GC-rich templates or betaine for AT-rich templates can reduce this bias (Baskaran et al., 1996; van Dijk, Jaszczyszyn, & Thermes, 2014; Kozarewa et al., 2009). Finally, the number of PCR cycles has also been shown to influence results: while a higher number of PCR cycles might increase the likelihood that rare molecules are observed, it could also skew abundance estimates by amplifying the biases described above (Casbon, Osborne, Brenner, & Lichtenstein, 2011; Weyrich et al., 2017), but this can vary (Kreihenwinkel et al., 2017; Vierna, Dona, Vizcaino, Serrano, & Jovani, 2017).

Here, we explore the potential of polymerase choice to influence the results of metabarcoding analyses, with particular reference to polymerase GC bias. We selected the *trnL* g/h primer set (Taberlet et al., 2007) as our universal barcoding primers for this evaluation, as the target *trnL* (P6 loop) locus of the chloroplast genome is commonly used for plant metabarcoding studies (Pornon et al., 2016; Sønstebo et al., 2010; Valentini et al., 2009). In addition, amplicons derived from this primer set are within the range of 50 and 150 base

pairs (bp), which is suitable for degraded environmental DNA and also fully sequenceable using short-read sequencing technologies. We performed metabarcoding on DNA extracted from soil collected from St. Paul Island, Alaska, and on mixtures of synthetic oligonucleotides whose inserts varied by GC content, using six polymerases, including those commonly used in metabarcoding. Using these experiments, we asked three questions: (i) Does polymerase GC preference affect relative abundance estimates in metabarcoding data? (ii) Are some polymerases more appropriate for metabarcoding-derived estimates of relative abundance than others? And (iii) Does GC bias affect occurrence estimates in metabarcoding experiments?

2 | MATERIALS AND METHODS

2.1 | Experimental design overview

We designed our experiment to ask three questions. First, *Does polymerase GC preference affect relative abundance estimates in metabarcoding data?* To answer this, we performed metabarcoding analyses of sedimentary DNA samples collected from St. Paul Island, Alaska. We performed two separate tests. First, we performed *trnL* (P6 loop) metabarcoding from nine samples and compared DNA-derived biodiversity estimates and biodiversity estimates based on above-ground survey data from the same sites. Next, for four of these nine sedimentary DNA samples, we explored whether relative abundance changed during the course of PCR amplification, following the design depicted in Figure 1. In both of these tests, we found that polymerase GC preference did affect relative abundance estimated. Our second question was therefore *Are some polymerases more appropriate for metabarcoding-derived estimates of relative abundance than others?* To answer this question, we amplified pools of synthetic oligonucleotides with a range of GC contents using six different polymerases and measured the precision with which each polymerase reconstructed the starting concentrations of each oligonucleotide pool. Our third question was *Does GC bias also affect occurrence estimates in metabarcoding experiments?* To answer this question, we again used the sedimentary DNA samples from St. Paul Island, Alaska, but this time performed metabarcoding using the polymerase identified in Question 2 as the least biased. We estimated the reproducibility of occurrence data using rarefaction analysis of ten replicate PCRs per sample.

2.2 | Data generation

2.2.1 | Environmental DNA from St Paul Island, Alaska

We collected soil samples from St. Paul Island, Alaska. This small (~114 km²), isolated island is situated ~450 km west of the coast of Alaska in the Bering Sea (~50.2°N, 170.2°W). St. Paul is the largest and most northerly island of the Pribilof Islands (Mungoven, 2005), has a low diversity of plants and terrestrial mammals (Colinvaux, 1981; Preble & McAtee, 1923) and completely lacks trees. We selected nine sampling sites that were spatially separate from each other, geologically distinct and appeared to be colonized by different vegetative communities. At each site, a 1 × 1 m quadrat was chosen. We removed a ~15 × 15 × 10 cm (L × W × D) volume of surface soil from the centre of each quadrat using a knife and trowel that we cleaned with ethanol between uses. We transferred ~10–20 g of soil to a sterile 50-ml falcon tube for eDNA analyses.

In addition to collecting sediment, we performed surveys of above-ground vegetation. We photographed the surface vegetation in each quadrat and performed a census of each taxon growing within the unit. We counted stems from each representative of each plant taxon and tallied the total for each unit (no counts exceeded 50). For very widespread and ubiquitous taxa, including spreading mat-forming types (e.g., mosses growing at the ground surface) and oversized plants with wide crowns, we estimated relative abundance based on percentage coverage within the unit. We identified the majority of common taxa in the field by comparison with a local collection curated at the St. Paul Public School and verified taxonomic assignments using Hultén's floras (Hultén, 1960, 1968). We collected representative samples of distinct or unknown taxa for later taxonomic verification, which we carried out using the relevant published floras along with online keys and floristics data (Hultén, 1960; Mungoven, 2005; Stotler & Crandall-Stotler, 2005; Talbot & Talbot, 1994; Walker et al., 2005). We converted the count data and the proportion of ground covered as a rank order (1 = 1–20% cover or <10 count; 2 = 21–40% or 10–24 count; 3 = 41–60% or 25–50 count; 4 = 61–80%; 5 = 81–100%) as a proxy for plant abundance at each sampling location.

We extracted environmental DNA from all nine soil samples using the MoBio PowerSoil DNA Isolation Kit (now called Qiagen DNeasy PowerSoil Kit), following the manufacturer's instructions. To

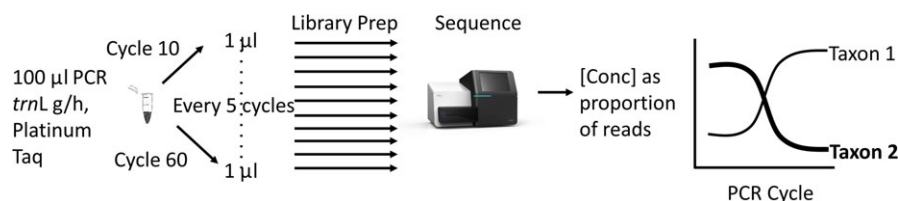


FIGURE 1 Schematic of the amplicon competition experiment. We chose four eDNA extracts and ran each in a PCR with *trnL* g/h and Platinum Taq using the recipe in Graham et al. (2016). Starting at cycle 10 and every five cycles up to cycle 60, we cooled the reaction to 20°C and removed 1 µl. We converted each 1 µl of PCR product into a sequenceable library individually. After sequencing and processing the reads, we plotted each amplicon as a function of PCR cycle and relative abundance

avoid contamination, we performed all steps in a clean laboratory that is physically isolated from other molecular biology research, while wearing sterile suits, face masks and gloves for DNA extractions and PCR set-up. To monitor cross-contamination, we extracted and processed the samples alongside two negative extraction controls, but did not use a positive control.

2.2.2 | Synthetic oligonucleotide pools

We designed and synthesized 12 oligonucleotides with inserts of 47 base pairs (bp) flanked by the *trnL* g/h primer binding sites with no mismatches (total length: 83 bp; Supplementary Table S1). This set included two oligonucleotides with 13% average GC content, two with 26% average GC content, two with 51% average GC content, two with 63% average GC content and four oligonucleotides with 38% average GC content. We then created six mixtures of these 12 oligonucleotides in which each oligonucleotide was included at different, but known, concentrations. We then diluted each mixture to 10 fM, which qPCR indicated was similar to the concentrations in our eDNA extracts. To verify pooling accuracy, we amplified each mixture using an approach that adds unique molecular identifiers (MIDs) to each starting molecule (Cole, Volden, Dharmadhikari, Scelfo-Dalbey, & Vollmers, 2016; Hoshino & Inagaki, 2017). Briefly, we first performed two cycles of PCR using modified versions of the *trnL* g/h primers that contained a 5' molecular identifier (which comprised five random nucleotides, followed by AT, followed by another three random nucleotides: NNNNNATNNN) and the Nextera adapter sequence (Supplementary Figure S1). This two-cycle PCR, which is performed using the permissive Phusion polymerase (New England Biosystems), adds to each starting molecule a uniquely identifying barcode that can be used to reconstruct bioinformatically the true starting relative abundance of molecules. After a clean-up step, we then amplified the product of this two-cycle PCR for an additional 30 cycles with standard Nextera indexing primers and the higher fidelity polymerase in Kapa HiFi ReadyMix (Kapa Biosystems). After sequencing, we counted the number of unique MIDs for each amplicon to verify the starting relative abundance of molecules in the pool.

2.2.3 | PCR amplification, library preparation, sequencing and bioinformatics

We performed PCR using the *trnL* g/h primers and six different polymerases (Table 1). We performed gradient PCR as necessary to determine optimal annealing temperatures for each of the different polymerases. For Platinum HiFi Taq, AmpliTaq Gold and Phusion, we used reagent mixes that are described in previous publications (Cole et al., 2016; De Barba et al., 2014; Graham et al., 2016). All final recipes and cycling conditions are provided in the supplement (Supplementary Table S2). We confirmed that amplification products were in the expected size range (50–150 bp) via gel electrophoresis, which also confirmed that all extraction and PCR-negative controls lacked visible amplification

products. We purified amplification products using a SPRI bead protocol (Rohland & Reich, 2012).

We transformed PCR amplicons into sequenceable libraries using two different approaches. Initially (for questions one and two), we used a lengthy protocol described by Meyer and Kircher (2010) (MK) that involves blunt-end repair, phosphorylation, adapter ligation and fill-in, and indexing PCR. To answer question three, we compared the MK protocol to a shorter and less expensive approach that amplifies DNA using *trnL* g/h primers with 5' overhangs containing the Illumina TruSeq adapter sequences. This made it possible to proceed directly to indexing PCR following the initial metabarcoding PCR, allowing library preparation to be completed in two steps (two PCR set-ups). To assess whether the two-step protocol performed differently from the MK protocol, we performed a comparative experiment in which we amplified DNA and sequenced libraries generated from a common master mix of Qiagen Multiplex Master Mix, water and template (consisting of an equimolar mixture of synthetic oligonucleotides). After sequencing, we found there was no significant difference between the two methods (standard least squares test: whole model F ratio = 0.55, p = 0.58, Supplementary Figure S2). While we find no difference between these two library preparation approaches, additional comparative analyses of prepared libraries for example using different GC content binning strategies, will be necessary to explore fully whether one library preparation approach is superior by all metrics to another.

For all experiments, we sequenced libraries on the Illumina MiSeq platform using 2×75 v3 chemistry, targeting 150,000 reads per sample. We used rarefaction to confirm that sequencing depth was sufficient to recover all amplified molecules (Hsieh et al., 2016).

After sequencing, we processed each data set using an in-house bioinformatics pipeline. Briefly, we removed adapters and merged overlapping reads using SEQPREP version 2 (<https://github.com/jstjohn/SeqPrep>), with the following flags: minimum length of reads (-L) 37 (combined length of the primer sequences plus one), overlap required to merge read1 and read2 (-o) 10, minimum length of adapter to consider trimming (-O) 8 and quality threshold (-q) 15. We filtered the merged reads and retained sequences containing either an exact match to the forward primer and the reverse complement of the reverse primer (correct orientation) or an exact match to the

TABLE 1 The six polymerases used in this study. *Platinum HiFi is a blend of two polymerases (one proofreading, one not)

Polymerase/mix	Manufacturer	Proofreading	Hot start
AmpliTaq Gold, Buffer II	Applied Biosystems	N	Y
Kapa HiFi ReadyMix	Kapa Biosystems	Y	Y
Phusion	New England BioLabs	Y	N
Platinum HiFi	Invitrogen	Y*	Y
Q5 2x Master Mix	New England BioLabs	Y	Y
Qiagen Multiplex Master Mix	Qiagen	N	Y

reverse primer and the reverse complement of the forward primer (incorrect orientation). We then reverse-complemented the data in the incorrect orientation using the `FASTX` toolkit version 0.0.13 (http://hannonlab.cshl.edu/fastx_toolkit/) and concatenated these data with those in the correct orientation. We trimmed any remaining adapter and PCR primer sequences from the ends of the filtered reads and removed any reads that retained any primer sequences or that were shorter than six base pairs using `PRINSEQ-lite` version 0.20.4 (Schmieder & Edwards, 2011). We created a single file with all unmerged reads, so that read1 and read2 were on the same line and processed this file as described above. We then split this file back into read1 and read2 files. We did not remove sequences that contained a mismatch to known synthetic oligonucleotide insert sequences (see below). For the sequence data derived from the St. Paul soil samples, all amplicons were short enough that the sequences could be merged. We used the `OBITOOLS` software defaults (Boyer et al., 2016) to group identical sequences (`obiuniq`), remove singletons and PCR artefacts (`obiclean -H`) and compare the sequences to the arctic, boreal and *embl* reference libraries (Sønsteby et al., 2010; Willerslev et al., 2014) to identify the reads to their best-associated plant taxa. Because we used three reference libraries, three separate result files were created for each sample (one for each reference library). We parsed the three files using a script that compared the results in each file and extracted only the entries with the highest percentage identity and lowest taxonomic rank. If two species of the same genus were seen, that sequence was classified to the genus level. We set a cut-off value of 98% identity and removed reads at proportions less than 0.001. The number of raw and merged reads and number of identified taxa per sample are listed in the Supplementary Materials (Supplementary Table S3).

For the synthetic oligonucleotide pools, we used `grep` to pull out the known sequences and their reverse complements and count how many times they occurred within each `FASTA` file. As the `OBITOOLS` and `GREP` methods both provided count data, we converted these counts to relative abundances.

2.3 | Data analysis

2.3.1 | Question 1: Does polymerase GC preference affect relative abundance estimates in metabarcoding data?

For the nine St. Paul samples, we performed ten replicate PCRs per sample using Platinum HiFi Taq polymerase (Invitrogen) following the protocol found in Graham et al. (2016). After sequencing and read processing as detailed above, we used standard least squares to test the effects of above-ground vegetation abundance and amplicon average GC content on DNA relative abundance, both separately and interactively.

To test the effect of PCR cycle number on the relative abundances of different plant taxa, we chose four St. Paul soil eDNA extracts and two PCR controls, scaled up the PCR to 100 μ L and collected 1 μ L aliquots at five-cycle intervals from cycles 10–60 (Figure 1). We used a large reaction volume to minimize the impact of aliquot removal and cooled the reaction to 20 C during each

collection step to avoid evaporation. Large numbers of cycles are often used in metabarcoding experiments because the target loci are at very low abundances relative to the total amount of extracted DNA and eDNA extracts often have PCR inhibitors (Kennedy, Callahan, & Carlson, 2013). We used 60 cycles to be sure that all PCRs had reached the plateau phase. Each aliquot was made into an Illumina sequencing library individually using a library preparation protocol based on Meyer and Kircher (2010) (as detailed above). We called this our amplicon competition experiment (Figure 1).

2.3.2 | Question 2: Are some polymerases more appropriate for metabarcoding-derived estimates of relative abundance than others?

We assessed whether six polymerases (Table 1) could individually maintain the starting ratio (relative abundance) of oligonucleotides in mixtures after 35 cycles of PCR. For each polymerase, we performed six experiments in which synthetic oligonucleotides were combined at different ratios based on sequence GC content. The oligonucleotides were combined (1) in equimolar ratios (two experiments), (2) by increasing proportion with GC content, (3) by decreasing proportion with GC content, (4) with extreme GC contents being most abundant and (5) with extreme GC contents being least abundant. For each experiment, we performed metabarcoding PCRs in triplicate using the *trnL g/h* primers. After obtaining relative abundance estimates for each oligonucleotide in each pool, we plotted expected abundances (relative abundance prior to amplification) versus observed abundances (relative abundance after amplification) for each polymerase. We then calculated the Pearson correlation coefficient between observed and expected abundance values for each enzyme.

2.3.3 | Question 3: Does GC bias affect occurrence estimates in metabarcoding experiments?

We again performed metabarcoding on the nine St. Paul soil eDNA extracts as described for Question 1, but used the Qiagen Multiplex Master Mix (Qiagen), which our results indicated is the least biased of the six polymerases tested (see below). As with the experiment described in Question 1 using Platinum HiFi Taq (Invitrogen), we performed ten replicate PCRs for each sample. We assigned amplicons to taxa as described above. We then performed rarefaction for each replicate set from both polymerases using *iNEXT* (Hsieh et al., 2016) in R version 3.4.2 (<http://www.R-project.org/>).

3 | RESULTS

3.1 | Question 1: Does polymerase preference for certain GC contents affect relative abundance estimates in metabarcoding data?

For this question, we used the above-ground vegetation abundance data, which were collected prior to the DNA work, and the Platinum HiFi Taq-amplified metabarcoding data. Both data sets were

generated from the same nine localities on St. Paul. Using both of these, we plotted all plant taxa that were identified using both above-ground and eDNA at all locations on the same plot but split into GC content bins (Figure 2). The x-values are above-ground ranked abundances, and the y-values are mean eDNA abundance across replicates. When we compared relative abundance estimates from the metabarcoding experiments to the relative abundance inferred from above-ground biomass, we found that whether or not these two estimates agreed depended on average GC content of the plant's *trnL* (P6 loop) locus (standard least squares, whole model: $F = 34.25$, $p < 0.0001$; effect tests: average GC, $t = 1.54$, $p = 0.124$, above-ground abundance, $t = 12.27$, $p < 0.0001$, average GC*above-ground abundance, $t = 4.39$, $p < 0.0001$). Figure 2 shows that above-ground and eDNA-based estimates of abundance are correlated most strongly in middle GC content bins, but this relationship decreases or disappears completely in the more extreme GC content bins. This pattern is consistent with the previously reported optimal GC content of 34–38% for Platinum HiFi Taq polymerase (Dabney & Meyer, 2012).

While this pattern observed in Figure 2 supports the hypothesis that sequences with certain GC contents are preferentially amplified via PCR, it does not exclude the possibility that biological factors, such as differences in above- versus below-ground biomass, are influencing the results. We therefore performed an additional experiment in which we measured changes in DNA-based relative abundance estimates directly during the course of PCR for four St. Paul eDNA extracts (Figure 1). Figure 3a shows the changes in relative abundance of the twelve most abundant taxa in each of the four samples during cycles 20 through 60 of the PCR. Libraries from cycles 10 and 15 had no sequenceable amplicon molecules. Exponential amplification appears to start at cycle 30 for all samples, and this was confirmed by qPCR (Supplementary Figure S3). We calculated the fold change from cycle 30–60 and used this to quantify the increase or decrease in the relative abundance of each amplicon. We then recorded the number of primer mismatches and barcode length for each amplicon. We found that neither primer mismatches nor amplicon length explained the increase or decrease in relative abundance (primer mismatches, $R^2: 0.011$; sequence length,

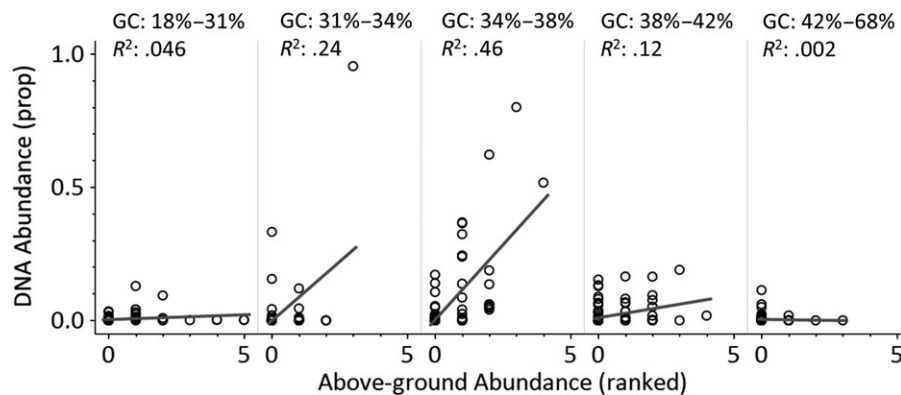


FIGURE 2 DNA abundance and above-ground abundance across average GC content bins. After collecting the DNA data, we took all data on plant taxa that were identified at all locations, put them in the same plot and split them into GC content bins. Each point is a plant taxon where x is its ranked above-ground abundance and y is its mean DNA abundance across replicates. Some taxa identified in the above-ground were not found in the DNA data and some found in the DNA were not found in the above-ground data. For above-ground abundance, 5 is the highest rank, meaning the most abundant, whereas 0 indicates absence. Lines are linear best fits with p -values > 0.3 for all bins except the middle bin where the $p = 0.03$

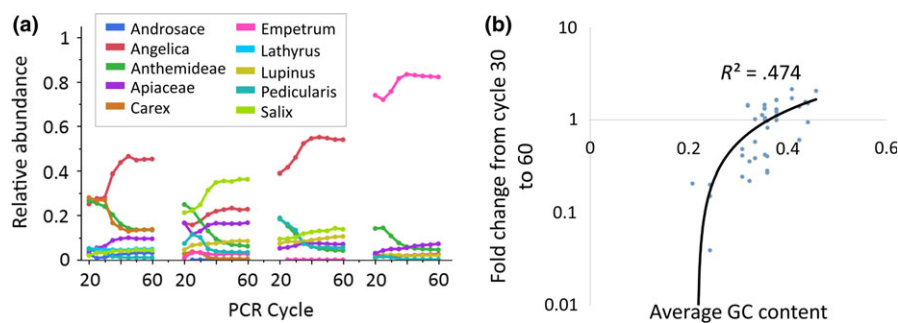


FIGURE 3 Changes in relative abundance over the duration of a 60-cycle PCR for four St. Paul Island sediment samples. (a) Plots showing relative abundance measured at 5-cycle intervals between cycles 20 and 60. Coloured lines show relative abundance estimates for the 10 most abundant plant taxa in these samples. (b) Plot describing the fold change for each taxon in each experiment between cycle 30 and cycle 60, with a linear line of best fit ($p = 0.002$), showing that change in relative abundance correlates with GC content. The y-axis is plotted on a log scale; therefore, values above 1 indicate that the amplicon is increasing in abundance from cycle 30–60 and values below 1 indicate that the amplicon is decreasing in abundance

R^2 : 0.095). However, we found a positive correlation with average GC content and fold change from cycle 30–60 (R^2 : 0.474, Linear fit p = 0.002; Figure 3b).

3.2 | Question 2: Are some polymerases more appropriate for metabarcoding-derived estimates of relative abundance than others?

Results from Question 1 suggest that Platinum HiFi Taq polymerase preferentially amplifies sequences with 34–38% GC. To identify polymerases that might be more appropriate for metabarcoding than Platinum HiFi Taq, we performed metabarcoding on mixtures of synthetic oligonucleotides with different GC contents using six commonly used polymerases (Table 1). We found that the correlation between observed and expected oligonucleotide proportions differed between enzymes (Figure 4). Among the polymerases tested, the Qiagen Multiplex Master Mix polymerase most accurately reconstructed the known starting relative abundances (Figure 4a, and varied the least in accuracy by GC content (Figure 4b). However, the Qiagen Multiplex Master Mix polymerase also had the highest proportion of sequences with at least one error (Figure 4c). Figure 5 shows the differences between observed and expected relative abundance using the most quantitatively accurate (Qiagen Multiplex Master Mix polymerase) and least quantitatively accurate (Phusion polymerase) enzymes. Detailed plots for the other four enzymes are provided in the Supplementary Materials (Supplementary Figures S4–S7).

3.3 | Question 3: Does GC bias affect occurrence data?

The results above show that polymerase biases can influence eDNA-based estimates of relative abundance. To test whether polymerase bias may also influence the accuracy of occurrence estimates, we performed an additional experiment in which we PCR-amplified the *trnL* (P6 loop) locus from the same nine St. Paul eDNA extracts that were amplified for Question 1, however, this time using the best-performing enzyme as identified by the synthetic oligonucleotide experiment above, Qiagen Multiplex Master Mix. As with Platinum HiFi Taq polymerase, we performed 10 replicate PCRs for each of the nine eDNA samples, and used rarefaction to confirm that sequencing depth of each PCR library was sufficient to recover all amplified molecules (Hsieh et al., 2016). We then performed additional rarefaction analyses, this time asking whether additional PCR replicates were contributing significantly towards biodiversity estimates. We found that after 10 replicates, mean sample coverage (the probability that all rare taxa have been recovered) was not significantly different when using the Qiagen Multiplex Master Mix compared to Platinum HiFi Taq (t = -0.66, df = 15.76, p = 0.52; Figure 6). In addition, despite the fact that St. Paul has low plant diversity (Colinvaux, 1981; Preble & McAtee, 1923), only one site appears to have reached a rarefaction plateau, which would suggest that the majority of species present have been sequenced after 10 replicates.

However, when we compared this to the data generated using Platinum HiFi Taq, this sample had not yet reached a rarefaction plateau. Given the small sample size, it is not possible to know whether this difference is due to polymerase choice or to chance.

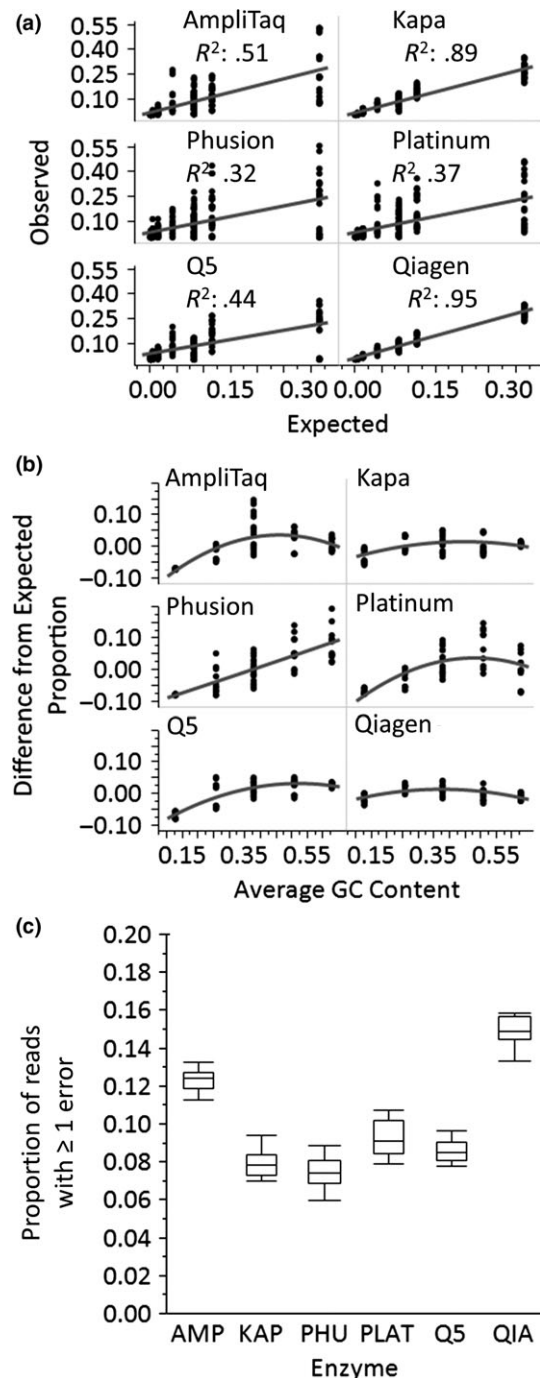


FIGURE 4 Testing polymerases using pools of synthetic oligonucleotides. a, b and c combine data from six pools of synthetic oligonucleotides amplified using six polymerases. (a) Observed proportions plotted against expected proportions for six polymerases. Each panel contains data for all six pools of oligos. (b) Difference from expected proportions plotted against average GC content. Here, we only used data from the equimolar pools. (c) Proportion of reads with at least one error for each enzyme/mix

4 | DISCUSSION

Our results show polymerase GC bias can dramatically alter the relative abundance of molecules during PCR. It is important, therefore, to use an experimental approach in metabarcoding that limits the influence of polymerase GC bias. Molecular identifier (MID), also called unique molecular identifier (UMI), methods (Cole et al., 2016) offer a possible solution, as they allow each starting molecule to be disambiguated bioinformatically after PCR. In this way, GC bias that manifests during PCR can be effectively ignored. However, these methods are not yet optimized for the mixed, low concentration samples that are most often available for metabarcoding. While we successfully tested a UMI approach for the analysis of synthetic mixtures of oligonucleotides, the approach often failed to produce sequencing libraries when analysing actual eDNA samples. This may be due to inhibitors and/or very low concentrations of target DNA

compared to all extracted DNA. Because polymerases vary in the degree to which they are biased towards GC content, another approach is to simply choose the least biased polymerase. Of the six polymerases evaluated here, our data show that the Qiagen Multiplex Master Mix is the least biased and effectively retains abundance ratios throughout the PCR (R^2 : 0.95). Qiagen Multiplex Master Mix (but not the enzyme, HotStarTaq, itself) was originally engineered for experiments that targeted multiple templates simultaneously, which may explain why it performs well here (Qiagen 2013).

If a biased polymerase is used in metabarcoding, the DNA results may not reflect the true relative abundance of target taxa. For the plant *tmL* (P6 loop) locus, for example, GC content varies considerably among major plant growth forms (Figure 7). The GC content of forbs, or low-lying herbaceous flowering plants, falls mainly within the range preferred by most polymerases (Dabney & Meyer, 2012).

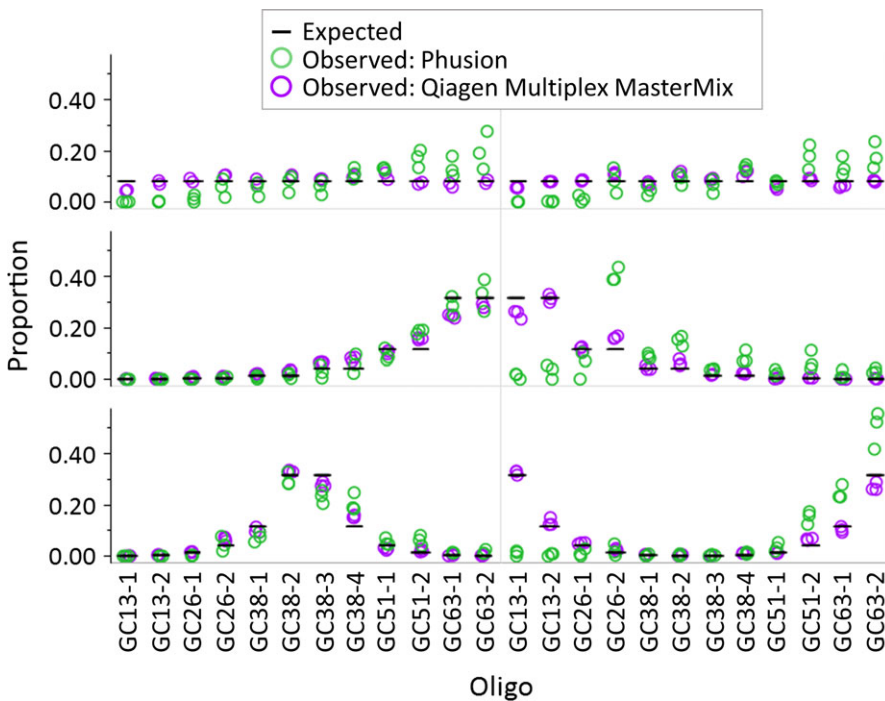


FIGURE 5 Expected (black lines) and observed abundances of the six synthetic oligonucleotide mixtures using Phusion (green open circles) and Qiagen Multiplex Master Mix (purple open circles) plotted as proportional data. Each oligonucleotide was pooled at 10 μ M, and then, each pool was diluted to 10 fM. Each 10 fM pool underwent PCR using the six different polymerases. The results for the best and the worst polymerases are plotted here

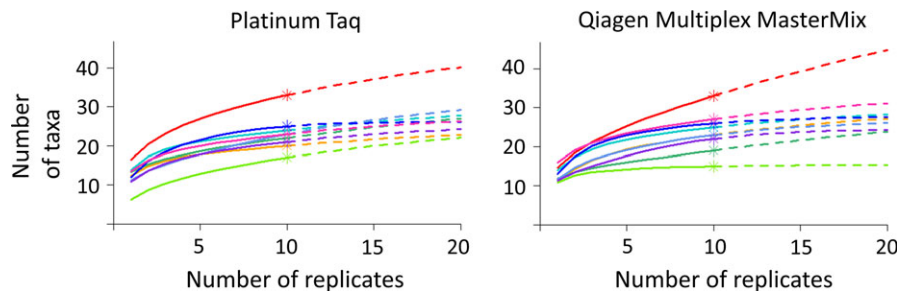


FIGURE 6 Rarefaction curves resulting from metabarcoding experiments for nine sites on St Paul Island, Alaska, using Platinum HiFi Taq as described in Graham et al. (2016) and the Qiagen Multiplex Master Mix following manufacturer's instructions. For each extract and polymerase, we performed 10 replicate PCRs. Rarefaction plots describe the number of unique taxa added per replicate. Solid lines are results from the 10 experiments, and dashed lines are predicted values calculated using iNEXT (Hsieh et al., 2016) in R version 3.4.2

Our DNA-based relative abundance estimates of plants from St. Paul (Figure 8) and those previously published from Siberia and Alaska (Supplementary Figure S8) (Willerslev et al., 2014) were both generated using Platinum HiFi Taq polymerase targeting the *trnL* P6-loop locus and showed that graminoids (grasses and sedges) were less abundant than forbs. Because this pattern falls within the biases of Platinum HiFi Taq polymerase, these results may simply reflect polymerase bias rather than true biological signal.

Although our results indicate that GC bias can confound metabarcoding-based relative abundance estimates, other potential sources of bias may also influence amplicon competition during PCR. For example, differences in the number of mismatches between the sequence and the primer at the primer binding site and differences in template length will also affect the efficiency with which an amplicon is copied (Stadhouders et al., 2010). While we did not find that the number of primer mismatches affected the efficiency of replication, few taxa have mismatches to the *trnL* g/h primers (Taberlet et al., 2007). Primer mismatches have been shown, however, to influence relative abundance for other metabarcoding loci (Piñol et al., 2015). In addition, shorter molecules tend to amplify more readily than longer molecules during PCR (Shagin, Lukyanov, Vagner, & Matz, 1999), and while most sequences amplified by the *trnL* g/h primers in this study tended to be around the same length, other metabarcoding loci vary considerably in barcode length between amplified taxa. Another source of bias during PCR is homopolymer repeats (Kieleczawa, 2006). In our amplicon competition experiment using Platinum HiFi Taq, the plant taxa Anthemideae and *Pedicularis* decreased in abundance in all four samples despite having optimal (Anthemideae has a GC content of 36%) and close to optimal (*Pedicularis* is 31%) GC contents, which may be because these barcodes contain 8- and 9-bp-long homopolymer runs, respectively. In comparison with Platinum HiFi Taq, we noted that Anthemideae and *Pedicularis* had increased abundances when using Qiagen Multiplex Master Mix (Supplementary Figures S9–S12), suggesting that Qiagen Multiplex Master Mix was not deterred by the homopolymer repeats. Finally, polymerase error rates are a potential source of error in metabarcoding experiments, and our results showed that HotStarTaq in the Qiagen Multiplex Master Mix had the highest error rate of the six polymerases used (Figure 4c). Polymerase error has the potential to produce false-positive results when barcoding loci differ by one or a few base pairs, although this may be ameliorated by bioinformatic pipelines capable of identifying potential sequencing errors.

Our results suggest that occurrence data, which have been believed to be largely reliable from metabarcoding experiments, can also be challenging to interpret. While it is understood that rare taxa may be more difficult to identify than common taxa, recommendations within the field have been to perform replicate PCRs, with little guidance as to how many PCRs are necessary. Our experiments from St. Paul suggest, however, that more than 10 replicate PCRs would be necessary to sample the breadth of taxa within our extracts, regardless of polymerase GC bias. In many instances, it may be more practical to combine DNA-based surveys with other data types, such

as pollen and identification of macroscopic remains (Birks & Birks, 2015). While site occupancy models offer a potential solution to estimate the number of replicates required to identify rare taxa (Dorazio & Erickson, 2017; Schmidt et al., 2013), these are constructed for single species and would not be practical for experiments that aim to describe an entire community. We note, however, that the most abundant taxa were recovered in all PCR replicates for all sites and both polymerases, suggesting that DNA metabarcoding is a reasonable approach to identify at least the most abundant taxa in an environment, even if only a single replicate PCR is performed (Leray & Knowlton, 2017).

While our current work has identified an experimental approach to reduce the influence of GC content on relative abundance estimates in metabarcoding, it is important also to consider other sources of potential biases and error when interpreting results. For example, errors such as tag switching, where sample-specific barcodes are associated with the incorrect sample during either library

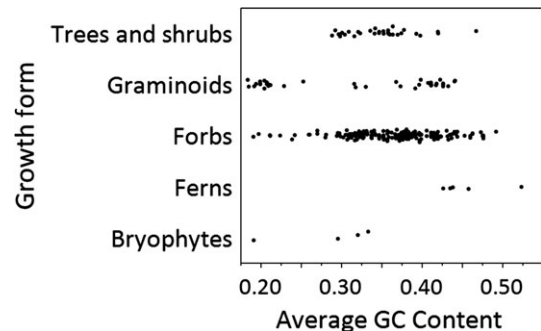


FIGURE 7 Average GC content across different plant growth forms. The data come from the current study and Willerslev et al. (2014). Both studies used Platinum HiFi Taq polymerase. Ferns include horsetails

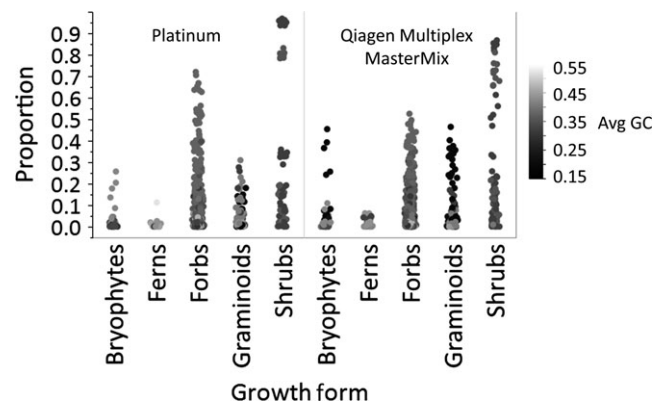


FIGURE 8 St Paul data generated using eDNA and separated into plant growth forms. Ferns include horsetails. All plant taxa from all samples are plotted both using Platinum HiFi Taq and Qiagen Multiplex Master Mix. Each data point is the relative abundance of a plant taxon from a particular location grouped into its growth form and shaded based on its average *trnL* p6-loop GC content. The darker the point is, the lower the average GC content

preparation (Schnell, Bohmann, & Gilbert, 2015) or sequencing (Kircher, Sawyer, & Meyer, 2012), may influence both occurrence and relative abundance data. Fortunately, the latter problem can be mitigated by adding indices to both ends of the molecule (Kircher et al., 2012). The choice of bioinformatic pipeline can also influence results. For example, in a recent analysis of the metagenome of fresh basil, three of four pipelines identified *Salmonella* but, because *Salmonella* was not identified via qPCR, the authors concluded that the bioinformatic results were erroneous (Ceuppens, De Coninck, Bottedoorn, Van Nieuwerburgh, & Uyttendaele, 2017). While public databases containing metabarcoding loci continue to expand in taxonomic depth (Bell, Loeffler, & Brosi, 2017), some lineages are more poorly represented. Finally, biological differences between species, including variation per cell or tissue type in the number of amplifiable loci (Morley & Nielsen, 2016), differences in organism size, seasonal senescence and behaviour, may all influence the probability that an organism will be represented in a particular environmental sample. Although work remains to be performed to better understand the consequences of these various types of bias and error, metabarcoding remains a powerful approach to quickly and inexpensively characterize communities.

5 | CONCLUSION

Despite the rapid growth of metabarcoding as a technique for characterizing communities from eDNA samples, relatively little attention has been given to validating the methodology and understanding its limitations. Polymerase GC bias is a known challenge for applications that rely on PCR (Aird et al., 2011; Dabney & Meyer, 2012; Kozarawa et al., 2009). With the advent of next-generation sequencing approaches, PCR-free methods have been developed to convert extracted DNA into sequenceable molecules (Kozarawa et al., 2009). PCR remains the most useful approach to catalogue diversity in environmental samples, however, as the number of target molecules is small relative to the total extracted DNA. For this reason, it is important to understand the influence of GC bias in metabarcoding approaches and, if possible, mitigate these biases. Here, we showed that many commonly used PCR protocols are not appropriate for generating reliable estimates of relative abundance. In these cases, our results show that the relative abundance of amplified sequences changes during PCR cycling and that these changes are related to the GC content of the target. Of the six polymerases and mixtures tested, Qiagen Multiplex Master Mix provided the most accurate estimates of relative abundance, but also generated the highest error rate. However, we found no evidence that occurrence data were influenced by polymerase bias.

ACKNOWLEDGEMENTS

We thank Soumaya Belmecheri for fieldwork assistance, Brendan O'Connell and Joshua Kapp for assistance with laboratory work and for writing scripts, and Duane Froese for discussion. We thank the

two reviewers whose comments helped improve this manuscript. This work was funded by awards from the Gordon and Betty Moore Foundation (GBMF-3804), the National Science Foundation (ARC-1203990) and the University of California Office of the President (20160713SC).

DATA ACCESSIBILITY

The processed sequence data and above-ground survey data are available in the Dryad repository (<https://doi.org/10.5061/dryad.k129c>). Raw amplicon sequence data have been made available on the NCBI Short Read Archive (BioProject: PRJNA433185).

AUTHOR CONTRIBUTIONS

B.S., P.D.H., L.A.N. and Y.W. designed and executed the collection of sediment samples. L.A.N. and Y.W. identified the plants. R.V.N., R.E.G., P.D.H. and B.S. designed the amplicon competition experiment and the synthetic oligonucleotides. C.V. and R.V.N. designed the experiments testing different polymerases. R.V.N. and M.L. performed the laboratory work. R.V.N. and C.V. analysed the data. R.V.N., B.S., C.V. and P.D.H. wrote the manuscript, with critical input from all remaining authors.

REFERENCES

- Aird, D., Ross, M. G., Chen, W.-S., Danielsson, M., Fennell, T., Russ, C., ... Gnrirke, A. (2011). Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries. *Genome Biology*, 12, R18.
- Andersen, K., Bird, K. L., Rasmussen, M., Haile, J., Breuning-Madsen, H., Kjaer, K. H., ... Willerslev, E. (2012). Meta-barcoding of "dirt" DNA from soil reflects vertebrate biodiversity. *Molecular Ecology*, 21, 1966–1979.
- Anderson-Carpenter, L. L., McLachlan, J. S., Jackson, S. T., Kuch, M., Lumibao, C. Y., & Poinar, H. N. (2011). Ancient DNA from lake sediments: Bridging the gap between paleoecology and genetics. *BMC Evolutionary Biology*, 11, 30.
- Barnes, M. A., & Turner, C. R. (2016). The ecology of environmental DNA and implications for conservation genetics. *Conservation Genetics*, 17(1), 1–17.
- Baskaran, N., Kandpal, R. P., Bhargava, A. K., Glynn, M. W., Bale, A., & Weissman, S. M. (1996). Uniform amplification of a mixture of deoxyribonucleic acids with varying GC content. *Genome Research*, 6, 633–638.
- Bell, K. L., Loeffler, V. M., & Brosi, B. J. (2017). An rbcL reference library to aid in the identification of plant species mixtures by DNA metabarcoding. *Applications in plant sciences*, 5, pii: apps.1600110.
- Birks, H. J. B., & Birks, H. H. (2015). How have studies of ancient DNA from sediments contributed to the reconstruction of Quaternary floras? *New Phytologist*, 209, 499–506.
- Boyer, F., Mercier, C., Bonin, A., Le Bras, Y., Taberlet, P., & Coissac, E. (2016). OBITOOLS: A unix-inspired software package for DNA metabarcoding. *Molecular Ecology Resources*, 16, 176–182.
- Buxton, A. S., Groombridge, J. J., Zakaria, N. B., & Griffiths, R. A. (2017). Seasonal variation in environmental DNA in relation to population size and environmental factors. *Scientific Reports*, 7, 46294.
- Casbon, J. A., Osborne, R. J., Brenner, S., & Lichtenstein, C. P. (2011). A method for counting PCR template molecules with application to next-generation sequencing. *Nucleic Acids Research*, 39, e81.

- Ceuppens, S., De Coninck, D., Botteldoorn, N., Van Nieuwerburgh, F., & Uyttendaele, M. (2017). Microbial community profiling of fresh basil and pitfalls in taxonomic assignment of enterobacterial pathogenic species based upon 16S rRNA amplicon sequencing. *International Journal of Food Microbiology*, 257, 148–156.
- Cole, C., Volden, R., Dharmadhikari, S., Scelfo-Dalbey, C., & Vollmers, C. (2016). Highly accurate sequencing of full-length immune repertoire amplicons using Tn5-enabled and molecular identifier-guided amplicon assembly. *Journal of Immunology*, 196, 2902–2907.
- Colinvaux, P. (1981). Historical ecology in beringia: The south land bridge coast at St. Paul Island. *Quaternary Research*, 16, 18–36.
- Cooper, A., & Poinar, H. N. (2000). Ancient DNA: Do it right or not at all. *Science*, 289, 1139.
- Corinaldesi, C., Barucca, M., Luna, G. M., & Dell'Anno, A. (2011). Preservation, origin and genetic imprint of extracellular DNA in permanently anoxic deep-sea sediments. *Molecular Ecology*, 20, 642–654.
- Dabney, J., & Meyer, M. (2012). Length and GC-biases during sequencing library amplification: a comparison of various polymerase-buffer systems with ancient and modern DNA sequencing libraries. *BioTechniques*, 52, 87–94.
- Darling, J. A., & Mahon, A. R. (2011). From molecules to management: Adopting DNA-based methods for monitoring biological invasions in aquatic environments. *Environmental Research*, 111, 978–988.
- De Barba, M., Miquel, C., Boyer, F., Mercier, C., Rioux, D., Coissac, E., & Taberlet, P. (2014). DNA metabarcoding multiplexing and validation of data accuracy for diet assessment: Application to omnivorous diet. *Molecular Ecology Resources*, 14, 306–323.
- Deagle, B. E., Gales, N. J., Evans, K., Jarman, S. N., Robinson, S., Trebilco, R., & Hindell, M. A. (2007). Studying seabird diet through genetic analysis of faeces: A case study on macaroni penguins (*Eudyptes chrysolophus*). *PLoS ONE*, 2, e831.
- Deagle, B. E., Thomas, A. C., Shaffer, A. K., Trites, A. W., & Jarman, S. N. (2013). Quantifying sequence proportions in a DNA-based diet study using Ion Torrent amplicon sequencing: Which counts count? *Molecular Ecology Resources*, 13, 620–633.
- Deiner, K., Bik, H. M., Mächler, E., Seymour, M., Lacoursière-Roussel, A., Altermatt, F., ... Bernatchez, L. (2017). Environmental DNA metabarcoding: Transforming how we survey animal and plant communities. *Molecular Ecology*, 26, 5872–5895.
- van Dijk, E. L., Jaszczyszyn, Y., & Thermes, C. (2014). Library preparation methods for next-generation sequencing: Tone down the bias. *Experimental Cell Research*, 322, 12–20.
- Direito, S. O. L., Marees, A., & Röling, W. F. M. (2012). Sensitive life detection strategies for low-biomass environments: Optimizing extraction of nucleic acids adsorbing to terrestrial and Mars analogue minerals. *FEMS Microbiology Ecology*, 81, 111–123.
- Dorazio, R. M., & Erickson, R. A. (2017). eDNA occupancy: An R package for multi-scale occupancy modeling of environmental DNA data. *Molecular Ecology Resources*, 18, 368–380.
- Dunn, N., Priestley, V., Herraiz, A., Arnold, R., & Savolainen, V. (2017). Behavior and season affect crayfish detection and density inference using environmental DNA. *Ecology and Evolution*, 7, 7777–7785.
- Elbrecht, V., & Leese, F. (2015). Can DNA-based ecosystem assessments quantify species abundance? Testing primer bias and biomass–sequence relationships with an innovative metabarcoding protocol. *PLoS ONE*, 10, 7. e0130324-16
- Fahner, N. A., Shokralla, S., Baird, D. J., & Hajibabaei, M. (2016). Large-scale monitoring of plants through environmental DNA metabarcoding of soil: recovery, resolution, and annotation of four DNA markers. *PLoS ONE*, 11, e0157505.
- Ficetola, G. F., Pansu, J., Bonin, A., Coissac, E., Giguët-Covex, C., De Barba, M., ... Taberlet, P. (2015). Replication levels, false presences and the estimation of the presence/absence from eDNA metabarcoding data. *Molecular Ecology Resources*, 15, 543–556.
- Giguët-Covex, C., Pansu, J., Arnaud, F., Rey, P. J., Griggo, C., Gielly, L., ... Taberlet, P. (2014). Long livestock farming history and human landscape shaping revealed by lake sediment DNA. *Nature Communications*, 5, 3211.
- Graham, R. W., Belmecheri, S., Choy, K., Culleton, B. J., Davies, L. J., Froese, D., ... Wooller, M. J. (2016). Timing and causes of mid-Holocene mammoth extinction on St. Paul Island, Alaska. *Proceedings of the National Academy of Sciences of the United States of America*, 113, 9310–9314.
- Haile, J., Froese, D. G., Macphee, R. D. E., Roberts, R. G., Arnold, L. J., Reyes, A. V., ... Willerslev, E. (2009). Ancient DNA reveals late survival of mammoth and horse in interior Alaska. *Proceedings of the National Academy of Sciences of the United States of America*, 106, 22352–22357.
- Haile, J., Holdaway, R., Oliver, K., Bunce, M., Gilbert, M. T., Nielsen, R., ... Willerslev, E. (2007). Ancient DNA chronology within sediment deposits: Are paleobiological reconstructions possible and is DNA leaching a factor? *Molecular Biology and Evolution*, 24, 982–989.
- Hemery, L. G., Politano, K. K., & Henkel, S. K. (2017). Assessing differences in macrofaunal assemblages as a factor of sieve mesh size, distance between samples, and time of sampling. *Environmental Monitoring and Assessment*, 189, 413.
- Hoshino, T., & Inagaki, F. (2017). Application of stochastic labeling with random-sequence barcodes for simultaneous quantification and sequencing of environmental 16S rRNA genes. *PLoS ONE*, 12(1), e0169431.
- Hsieh, T. C., Ma, K. H., & Chao, A. (2016). iNEXT: An R package for rarefaction and extrapolation of species diversity (Hill numbers). *Methods in Ecology and Evolution/British Ecological Society*, 7, 1451–1456.
- Hultén, E. (1960). *Flora of the Aleutian Islands and westernmost Alaska Peninsula: With notes on the Flora of commander Islands*. New York, NY: Hafner Pub. Co..
- Hultén, E. (1968). *Flora of Alaska and neighboring territories: A manual of the vascular plants*. Stanford, CA: Stanford University Press.
- Jerde, C. L., Mahon, A. R., Chadderton, W. L., & Lodge, D. M. (2011). "Sight-unseen" detection of rare aquatic species using environmental DNA. *Conservation Letters*, 4, 150–157.
- Jørgensen, T., Kjaer, K. H., Haile, J., Rasmussen, M., Boessenkool, S., Andersen, K., ... Willerslev, E. (2012). Islands in the ice: Detecting past vegetation on Greenlandic nunataks using historical records and sedimentary ancient DNA meta-barcoding. *Molecular Ecology*, 21, 1980–1988.
- Kennedy, S., Callahan, H., & Carlson, M. (2013). *Tips and tricks for isolation of DNA & RNA from challenging samples*. Carlsbad, CA: MO BIO Laboratories Inc.
- Kieleczawa, J. (2006). Fundamentals of sequencing of difficult templates—an overview. *Journal of Biomolecular Techniques*, 17, 207–217.
- Kircher, M., Sawyer, S., & Meyer, M. (2012). Double indexing overcomes inaccuracies in multiplex sequencing on the Illumina platform. *Nucleic Acids Research*, 40, e3.
- Kowalczyk, R., Taberlet, P., Kaminski, T., & Wojcik, J. (2011). Influence of management practices on large herbivore diet—Case of European bison in Białowieża Primeval Forest (Poland). *Forest Ecology and Management*, 261, 821–828.
- Kozarewa, I., Ning, Z., Quail, M. A., Sanders, M. J., Berriman, M., & Turner, D. J. (2009). Amplification-free Illumina sequencing-library preparation facilitates improved mapping and assembly of (G+C)-biased genomes. *Nature Methods*, 6, 291–295.
- Krehsenwinkel, H., Wolf, M., Lim, J. Y., Rominger, A. J., Simison, W. B., & Gillespie, R. (2017). Estimating and mitigating amplification bias in qualitative and quantitative arthropod metabarcoding. *Scientific Reports*, 7, 1–12.
- Lahoz-Monfort, J. J., Guillera-Arroita, G., & Tingley, R. (2016). Statistical approaches to account for false-positive errors in environmental DNA samples. *Molecular Ecology Resources*, 16, 673–685.

- Leray, M., & Knowlton, N. (2017). Random sampling causes the low reproducibility of rare eukaryotic OTUs in Illumina COI metabarcoding. *PeerJ*, 5, e3006–e3027.
- Marcelino, V. R., & Verbruggen, H. (2016). Multi-marker metabarcoding of coral skeletons reveals a rich microbiome and diverse evolutionary origins of endolithic algae. *Scientific Reports*, 6, 31508.
- Meyer, M., & Kircher, M. (2010). Illumina sequencing library preparation for highly multiplexed target capture and sequencing. *Cold Spring Harbor Protocols*, 6, pdb.prot5448.
- Morley, S. A., & Nielsen, B. L. (2016). Chloroplast DNA copy number changes during plant development in organelle DNA polymerase mutants. *Frontiers in Plant Science*, 7, 57.
- Mungoven, M. (2005) *Soil Survey of Saint Paul Island Area, Alaska*. National Resources Conservation Service, United States Department of Agriculture.
- Niemeyer, B., Epp, L. S., Stof- Leichsenring, K. R., Pstryakova, L. A., & Herzsuh, U. (2017). A comparison of sedimentary DNA and pollen from lake sediments in recording vegetation composition at the Siberian treeline. *Molecular Ecology Resources*, 17, e46–e62.
- Pawluczyk, M., Weiss, J., Links, M. G., Egaña Aranguren, M., Wilkinson, M. D., & Egea-Cortines, M. (2015). Quantitative evaluation of bias in PCR amplification and next-generation sequencing derived from metabarcoding samples. *Analytical and Bioanalytical Chemistry*, 407, 1841–1848.
- Pedersen, M. W., Overballe-Petersen, S., Ermini, L., Der Sarkissian, C., Haile, J., Hellstrom, M., ... Schnell, I. B. (2015). Ancient and modern environmental DNA. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, 370, 20130383.
- Pedersen, M. W., Ruter, A., Schweger, C., Friebe, H., Staff, R. A., Kjeldsen, K. K., ... Willerslev, E. (2016). Postglacial viability and colonization in North America's ice-free corridor. *Nature*, 537, 45–49.
- Piñol, J., Mir, G., Gomez-Polo, P., & Agustí, N. (2015). Universal and blocking primer mismatches limit the use of high-throughput DNA sequencing for the quantitative metabarcoding of arthropods. *Molecular Ecology Resources*, 15, 819–830.
- Pinto, A. J., & Raskin, L. (2012). PCR biases distort bacterial and archaeal community structure in pyrosequencing datasets. *PLoS ONE*, 7, e43093.
- Polz, M. F., & Cavanaugh, C. M. (1998). Bias in template-to-product ratios in multitemplate PCR. *Applied and Environmental Microbiology*, 64, 3724–3730.
- Pornon, A., Escaravage, N., Burrus, M., Holota, H., Khimoun, A., Mariette, J., ... Andalo, C. (2016). Using metabarcoding to reveal and quantify plant-pollinator interactions. *Scientific Reports*, 6, 27282.
- Preble, E. A., & McAtee, W. L. (1923). *A Biological Survey of the Pribilof Islands, Alaska, North American Fauna, No. 46*. Department of Agriculture, Bureau of Biological Survey, Government Printing Office, Washington DC, USA.
- Qiagen (2013). Qiagen Multiplex PCR Kit: Product Details.
- Rawlence, N. J., Lowe, D. J., Wood, J. R., Young, J. M., Jock Churchman, G., Huang, Y. T., & Cooper, A. (2014). Using palaeoenvironmental DNA to reconstruct past environments: Progress and prospects. *Journal of Quaternary Science*, 29, 610–626.
- Rees, H. C., Baker, C. A., Gardner, D. S., Maddison, B. C., & Gough, K. C. (2017). The detection of great crested newts year round via environmental DNA analysis. *BMC Research Notes*, 10, 327.
- Rohland, N., & Reich, D. (2012). Cost-effective, high-throughput DNA sequencing libraries for multiplexed target capture. *Genome Research*, 22, 939–946.
- Sanders, H. L. (1968). Marine benthic diversity: A comparative study. *American Naturalist*, 102, 243–282.
- Schmidt, B. R., Kéry, M., Ursenbacher, S., Hyman, O. J., & Collins, J. P. (2013). Site occupancy models in the analysis of environmental DNA presence/absence surveys: A case study of an emerging amphibian pathogen. *Methods in Ecology and Evolution/British Ecological Society*, 4, 646–653.
- Schmieder, R., & Edwards, R. (2011). Quality control and preprocessing of metagenomic datasets. *Bioinformatics*, 27(6), 863–864.
- Schnell, I. B., Bohmann, K., & Gilbert, M. T. P. (2015). Tag jumps illuminated—reducing sequence-to-sample misidentifications in metabarcoding studies. *Molecular Ecology Resources*, 15, 1289–1303.
- Shagin, D. A., Lukyanov, K. A., Vagner, L. L., & Matz, M. V. (1999). Regulation of average length of complex PCR product. *Nucleic Acids Research*, 27, e23.
- Shaw, J. L. A., Clarke, L. J., Wedderburn, S. D., Barnes, T. C., Weyrich, L. S., & Cooper, A. (2016). Comparison of environmental DNA metabarcoding and conventional fish survey methods in a river system. *Biological Conservation*, 197, 131–138.
- Sigsgaard, E. E., Nielsen, I. B., Bach, S. S., Lorenzen, E. D., Robinson, D. P., Knudsen, S. W., ... Thomsen, P. F. (2016). Population characteristics of a large whale shark aggregation inferred from seawater environmental DNA. *Nature Ecology & Evolution*, 1, 0004.
- Soininen, E. M., Valentini, A., Coissac, E., Miquel, C., Gielly, L., Brochmann, C., ... Taberlet, P. (2009). Analysing diet of small herbivores: The efficiency of DNA barcoding coupled with high-throughput pyrosequencing for deciphering the composition of complex plant mixtures. *Frontiers in Zoology*, 6, 16.
- Sønsteby, J. H., Gielly, L., Brysting, A. K., Elven, R., Edwards, M., Haile, J., ... Brochmann, C. (2010). Using next-generation sequencing for molecular reconstruction of past Arctic vegetation and climate. *Molecular Ecology Resources*, 10, 1009–1018.
- Stadhouders, R., Pas, S. D., Anber, J., Voermans, J., Mes, T. H., & Schutten, M. (2010). The effect of primer-template mismatches on the detection and quantification of nucleic acids using the 5' nuclease assay. *Journal of Molecular Diagnostics*, 12, 109–117.
- Stotler, R. E., & Crandall-Stotler, B. J. (2005). Bryophytes. Department of Plant Biology, Southern Illinois University.
- Suzuki, M. T., & Giovannoni, S. J. (1996). Bias caused by template annealing in the amplification of mixtures of 16S rRNA genes by PCR. *Applied and Environmental Microbiology*, 62, 625–630.
- Taberlet, P., Coissac, E., Pompanon, F., Gielly, L., Miquel, C., Valentini, A., ... Willerslev, E. (2007). Power and limitations of the chloroplast trnL (UAA) intron for plant DNA barcoding. *Nucleic Acids Research*, 35, e14.
- Talbot, S. S., & Talbot, S. L. (1994). Numerical classification of the coastal vegetation of Attu Island, Aleutian Islands, Alaska. *Journal of Vegetation Science*, 5, 867–876.
- Ushio, M., Fukuda, H., Inoue, T., Makoto, K., Kishida, O., Sato, K., ... Miya, M. (2017). Environmental DNA enables detection of terrestrial mammals from forest pond water. *Molecular Ecology Resources*, 17, e63–e75.
- Valentini, A., Miquel, C., Nawaz, M. A., Bellemain, E., Coissac, E., Pompanon, F., ... Taberlet, P. (2009). New perspectives in diet analysis based on DNA barcoding and parallel pyrosequencing: The trnL approach. *Molecular Ecology Resources*, 9, 51–60.
- Vierna, J., Dona, J., Vizcaino, A., Serrano, D., & Jovani, R. (2017). PCR cycles above routine numbers do not compromise high-throughput DNA barcoding results. *Genome*, 60(10), 868–873.
- Walker, D. A., Raynolds, M. K., Daniëls, F. J. A., Einarsson, E., Elvebakk, A., Gould, W. A., ... Yurtsev, A. (2005). The circumpolar Arctic vegetation map. *Vegetation Science*, 5, 757–764.
- Weyrich, L. S., Duchene, S., Soubrier, J., Arriola, L., Llamas, B., Breen, J., ... Cooper, A. (2017). Neanderthal behaviour, diet, and disease inferred from ancient DNA in dental calculus. *Nature*, 544, 357–361.
- Willerslev, E., Cappellini, E., Boomsma, W., Nielsen, R., Hebsgaard, M. B., Brand, T. B., ... Collins, M. J. (2007). Ancient biomolecules from deep ice cores reveal a forested southern Greenland. *Science*, 317, 111–114.
- Willerslev, E., Davison, J., Moora, M., Zobel, M., Coissac, E., Edwards, M. E., ... Taberlet, P. (2014). Fifty thousand years of Arctic vegetation and megafaunal diet. *Nature*, 506, 47–51.

- Willerslev, E., Hansen, A. J., Binladen, J., Brand, T. B., Gilbert, M. T., Shapiro, B., ... Cooper, A. (2003). Diverse plant and animal genetic records from Holocene and Pleistocene sediments. *Science*, *300*, 791–795.
- Yoccoz, N. G., Bråthen, K. A., Gielly, L., Haile, J., Edwards, M. E., Goslar, T., ... Taberlet, P. (2012). DNA from soil mirrors plant taxonomic and growth form diversity. *Molecular Ecology*, *21*, 3647–3655.
- Zielińska, S., Radkowski, P., Blendowska, A., Ludwig-Gałęzowska, A., Łoś, J. M., & Łoś, M. (2017). The choice of the DNA extraction method may influence the outcome of the soil microbial community structure analysis. *Microbiology Open*, *6*, <https://doi.org/10.1002/mbo3.453>.

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

How to cite this article: Nichols RV, Vollmers C, Newsom LA, et al. Minimizing polymerase biases in metabarcoding. *Mol Ecol Resour.* 2018;00:1–13. <https://doi.org/10.1111/1755-0998.12895>