

# Disentangling bias for non-destructive insect metabarcoding

**Francesco Martoni**<sup>Corresp., 1</sup>, **Alexander M Piper**<sup>1, 2</sup>, **Brendan C Rodoni**<sup>1, 2</sup>, **Mark J Blacket**<sup>1</sup>

<sup>1</sup> Agriculture Victoria Research, AgriBio Centre for AgriBio Science, State Government Victoria, Bundoora, Victoria, Australia

<sup>2</sup> School of Applied Systems Biology, La Trobe University, Bundoora, Victoria, Australia

Corresponding Author: Francesco Martoni

Email address: francesco.martoni@ecodev.vic.gov.au

A fast and reliable method for obtaining a species-level identification is a fundamental requirement for a wide range of activities, from plant protection and invasive species management to biodiversity assessments and ecological studies. For insects, novel molecular techniques such as DNA metabarcoding have emerged as a rapid alternative to traditional morphological identification, reducing the dependence on limited taxonomic experts. Until recently, molecular techniques have required a destructive DNA extraction, precluding the possibility of preserving voucher specimens for future studies, or species descriptions. Here we paired insect metabarcoding with two recent non-destructive DNA extraction protocols, to obtain a rapid and high-throughput taxonomic identification of diverse insect taxa while retaining a physical voucher specimen. The aim of this work was to explore how non-destructive extraction protocols impact the semi-quantitative nature of metabarcoding, which alongside species presence/absence also provides a quantitative, but biased, representation of their relative abundances. By using a series of mock communities representing each stage of a typical metabarcoding workflow we were able to determine how different morphological (i.e., insect size and hardness) and molecular traits (i.e., DNA extraction and PCR), interact with different protocol steps to introduce quantitative bias into non-destructive metabarcoding results. We discuss the relevance of taxonomic bias to metabarcoding identification of insects and potential approaches to account for it.

# Disentangling bias for non-destructive insect metabarcoding.

Francesco Martoni<sup>1</sup>, Alexander M. Piper<sup>1,2</sup>, Brendan C. Rodoni<sup>1,2</sup>, Mark J. Blacket<sup>1</sup>

<sup>1</sup> Agriculture Victoria Research, AgriBio Centre for AgriBio Science, Bundoora 3083, Victoria, Australia

<sup>2</sup> School of Applied Systems Biology, La Trobe University, Bundoora 3083, Victoria, Australia

Corresponding Author:

Francesco Martoni

5 Ring Road, Bundoora, Victoria, 3083, Australia.

Email address: [Francesco.Martoni@agriculture.vic.gov.au](mailto:Francesco.Martoni@agriculture.vic.gov.au).

## Abstract

A fast and reliable method for obtaining a species-level identification is a fundamental requirement for a wide range of activities, from plant protection and invasive species management to biodiversity assessments and ecological studies. For insects, novel molecular techniques such as DNA metabarcoding have emerged as a rapid alternative to traditional morphological identification, reducing the dependence on limited taxonomic experts. Until recently, molecular techniques have required a destructive DNA extraction, precluding the possibility of preserving voucher specimens for future studies, or species descriptions. Here we paired insect metabarcoding with two recent non-destructive DNA extraction protocols, to obtain a rapid and high-throughput taxonomic identification of diverse insect taxa while retaining a physical voucher specimen. The aim of this work was to explore how non-destructive extraction protocols impact the semi-quantitative nature of metabarcoding, which alongside species presence/absence also provides a quantitative, but biased, representation of their relative abundances. By using a series of mock communities representing each stage of a typical metabarcoding workflow we were able to determine how different morphological (i.e., insect size and hardness) and molecular traits (i.e., DNA extraction and PCR), interact with different protocol steps to introduce quantitative bias into non-destructive metabarcoding results. We discuss the relevance of taxonomic bias to metabarcoding identification of insects and potential approaches to account for it.

# Introduction

Species identification is a fundamental pre-requisite for basic and applied ecology. In the field of entomology, species-level identification is required for biodiversity assessments and checklists (Girón & Short 2021), understanding ecology and behavior (Lefort *et al.* 2020), forensic investigation (Pohjoismäki *et al.* 2010), taxonomy (Schutze *et al.* 2017), and management of agricultural pests. Invasive insect species are becoming a major threat to agroecosystems (Paini *et al.* 2016), with biological invasions becoming one of the main menaces to agricultural production (Meyerson & Mooney 2007; Hulme 2009; Chown *et al.* 2014). Therefore, extensive trapping and monitoring activities to detect new insect invasions are carried out both on agricultural properties (Low-Choy 2015) and protected environments, such as National Parks (e.g., Davidovitch *et al.* 2009). This surveillance is leading to an increasing demand of species-level identification for large volumes of insects trapped for plant protection, biosecurity and agriculture (Piper *et al.* 2019).

However, availability of taxonomic expertise for insect identification is extremely limited, with often only a few experts worldwide per taxonomic group. Therefore, a range of molecular techniques have been developed to allow more standardized identification of insect species by non-specialists (Piper *et al.* 2019). Most notably, the DNA barcoding technique (Hebert *et al.* 2003) allows comparison of a short standardized genetic region from an unidentified specimen to a vast number of known species deposited in reference databases. Generally, invertebrate barcoding studies have targeted the subunit 1 of the mitochondrial cytochrome oxidase gene (COI) for specimen identification (e.g., Andújar *et al.* 2018; Elbrecht *et al.* 2017; Yu *et al.* 2012). DNA barcoding is regularly applied to identification of undescribed species native to localised areas (Martoni *et al.* 2020), as well as to identification of invasive insect pests (Armstrong & Ball 2005), and is now widely accepted within plant pest diagnostics protocols (EPPO 2021; Ashfaq *et al.* 2016). However, difficulties remain scaling this approach to the sheer numbers of specimens that can be caught in a surveillance trap, or in a field sampling season. Therefore, with the advent of high throughput sequencing technologies, the focus is now shifting from the single specimen sequencing of DNA barcoding to identifying entire communities of species in parallel using DNA metabarcoding (Piper *et al.* 2019 and references therein).

The semi-quantitative nature of metabarcoding data has led to concerns around the appropriateness of this technology for surveillance, biomonitoring and assessment of pest pressures (Darling *et al.* 2020; Martins *et al.* 2019). Metabarcoding assays only provide relative abundance data, with sequence reads returned for a species only meaningful relative to the rest of the taxa within the sample (Gloor *et al.* 2017). In addition to quantification issues relating to sample composition, taxonomic biases can be introduced during laboratory processing of samples, including during DNA extraction procedures and genetic marker isolation and amplification. These biases arise due to species-specific differences in morphological and molecular traits which interact with steps of the laboratory protocol to preferentially detect certain taxa at the expense of others. DNA extraction from complex insect communities for

metabarcoding analysis has often involved destructive homogenisation of tissues, which results in larger-sized organisms contributing a larger quantity of DNA molecules to the DNA extraction pool than smaller bodied insects (Elbrecht *et al.* 2017). Nevertheless, when individual specimen size is accounted for through size sorting prior to DNA extraction, the influence of primer-template mismatch generally outweighs DNA extraction bias for macroinvertebrates (Braukmann *et al.* 2019; Elbrecht & Leese 2015), at least for destructive homogenisation-based DNA extraction.

Recently, non-destructive DNA extraction has emerged as an alternative to homogenisation-based methods in order to retain voucher specimens for morphological confirmation of metabarcoding detections (Carew *et al.* 2018; Martins *et al.* 2019; Batovska *et al.* 2021). This is of particular importance in the context of agricultural biosecurity and other regulatory applications of metabarcoding, allowing DNA sequences to be linked to an insect sample, which can be preserved in entomological collections for future records (Martoni *et al.* 2021a). While it is well established for homogenised samples that large organisms will have a higher abundance of DNA than small organisms, this may differ for non-destructive metabarcoding, where contribution of DNA may instead depend on surface/volume ratio of insect to extraction buffer (Marquina *et al.* 2019). Furthermore, differences in sclerotization of exoskeletons could affect permeability of DNA and impact detection efficiency (Carew *et al.* 2018; Marquina *et al.* 2019). Moreover, with this move from destructive to non-destructive DNA extraction it is unclear if earlier results and assumptions about the bias generating process still hold.

Here we compared two different non-destructive DNA extraction methods previously tested for their capability to extract DNA from preserved trapped insects (Martoni *et al.* 2021b). We applied these to mock communities composed of a mix of possible insect pests and harmless by-catch species, following a typical metabarcoding protocol from the recent literature, and using two combinations of degenerate PCR primers. We aimed to measure the taxonomic bias introduced by the DNA extraction, PCR, and library preparation stages, and evaluate the downstream effects on the two main diagnostic-related aspects: sensitivity and quantitation.

## Materials & Methods

### Samples and morphological traits

For this study we used adult specimens belonging to 16 insect species (Table 1). Of these, 15 species were obtained from insect colonies reared at the AgriBio laboratory of Agriculture Victoria, while another species (*Acizzia* sp.) was field-collected (Table 1). Insect specimens from the colonies were preserved in absolute ethanol and deposited into the Victorian Agricultural Insect Collection (VAIC). Measurements to obtain the volume size were taken from ten individuals for each species and the result was averaged in Table 1. Measurements were taken using the Leica Application Suite software v4.5.0, from five to 20 stacked images obtained using

a Leica stereo microscope M205C with a DFC450 camera. For each insect, images were taken from dorsal, lateral, and frontal view in order to obtain a measurement of the volume (length \* width \* depth) of head, thorax and abdomen. Hardness of the exoskeleton was estimated by attempting to pierce it using forceps, and the 16 species were assigned to three discrete categories (soft, intermediate, hard). Specimens from all taxa were then grouped into eight insect pools, with six pools containing 100 individual insects and the same 15 species, while another two pools also contained a single specimen of the *Bradysia* nr. *ocellaris*, for a total of 16 species, and 101 individuals (Table SM1).

## Molecular analysis

### Preparation of Insect mock communities.

In order to partition total protocol bias into its constituent steps, mock communities of bulk adult insect, genomic DNA and PCR amplicons were assembled to simulate the input of each major metabarcoding laboratory step.

DNA was non-destructively extracted from the eight pools of insects using both the QuickExtract kit (Biosearch Technologies, California, USA), for pools 1-4, and the DNEasy Blood and Tissue kit (Qiagen, Germany; “DNEasy” hereafter), for pools 5-8, following the methods used for single insects in Martoni *et al.* (2021b). These pools, while having a different number of individuals for some of the species, were prepared with an almost identical Order-level composition, taking into account insect biomass and exoskeleton hardness. Non-destructive DNA extraction using QuickExtract was performed as follows: ethanol was removed from the pooled insects and air-dried in tubes for 10 minutes. Five hundred microlitres of QuickExtract were added to the pooled insects, ensuring all insects were submerged, vortexed for 30 seconds, incubated at 65 °C for 6 minutes, vortexed for 15 seconds and incubated at 98 °C for 2 minutes. The supernatant containing the extracted DNA was then transferred to a new Eppendorf tube and stored at -20 °C until PCR amplification.

Non-destructive DNA extraction using DNEasy was performed following the first steps of the protocol presented in Martoni *et al.* (2019) and Bahder *et al.* (2015). Briefly, ethanol was removed from the insect pools (as above), insects were then submerged in an ATL buffer/Proteinase K mix with a ratio of 9:1 and then incubated for approximately 17 hours (overnight) at 56 °C. The supernatant was then removed from the insects (as above), and processed further following the manufacturer instructions (i.e., filter column purification and elution). Finally, pooled insect specimens were resuspended in absolute ethanol and preserved.

PCR amplification was conducted using either fwhF2 (GGDACWGGWTGAACWGTWTAYCCHCC)- fwhR2n (GTRATWGCHCCDGCTARWACWGG) or fwhF2 – HexCOIR4 (TATDGTRATDGCHCCNGC), which amplify almost entirely overlapping regions of COI (Vamos *et al.* 2017; Marquina *et al.* 2018). The primer pair fwhF2-fwhR2n targets a 205 bp amplicon (excluding primers) from 346 bp to 551 bp within the conventional COI barcode region, while fwhF2-HexCOIR4 target a 214 bp amplicon (excluding primers) from 346 bp to

560 bp. PCR amplification was performed using the Bioline MyFi DNA Polymerase kit (Meridian Bioscience, Ohio, United States of America) using 2.5  $\mu$ L of DNA template and 1  $\mu$ L each for the two primers (10 mM) in a 25  $\mu$ L final volume. The PCR was run with the same cycling conditions for both primer pairs, with an initial 5-minute denaturation at 95 °C, followed by 30 cycles, of denaturation at 95 °C for 45 seconds, annealing at 50 °C for 30 seconds and extension at 72 °C for 30 seconds, followed by a final extension at 72 °C for 7 minutes. PCR amplification was verified on a 1% w/v agarose gel.

PCR amplicons were used as template for a second round of PCR to attach Illumina sequencing adapters with unique dual indexes to each sample, using Phusion High-Fidelity DNA Polymerase (New England Biolabs, Massachusetts, USA). PCR conditions were an initial denaturation of 30 seconds at 98 °C followed by 8 cycles of denaturation at 98 °C for 10 seconds, annealing at 65 °C for 30 seconds and elongation at 72 °C for 30 seconds. The adapter tailed and indexed amplicons were purified using AMPure XP beads (Beckman Coulter, California, USA) following the manufacturer instructions. Library fragment size (amplicon + adapters) and absence of primer dimers was verified on an Agilent TapeStation (Agilent Technologies, California, USA) and all libraries were equimolarly pooled based on their concentrations as determined by Qubit dsDNA HS Fluorometric Quantification (ThermoFisher Scientific, Massachusetts, USA). As DNA concentrations in negative controls were too low to be measured, they were pooled at the same volume of the lowest concentration mock community library. The final pooled library was then diluted to 7 pM, spiked with 15% PhiX, and sequenced on the Illumina MiSeq platform using the V2 reagent kit (2 x 250 bp reads) (Illumina, California, USA).

#### *Preparation of DNA and PCR mock communities.*

To assemble mock communities representative of the post-DNA extraction stage (Figure 1; “DNA Pools”), DNA for each of the 16 species was destructively extracted separately from 5-20 homogenised individuals (depending on their size) with the DNEasy kit. The DNA from each insect species was then quantified using a Qubit 2.0 Fluorometer (Thermo Fisher Scientific, Massachusetts, USA) and pooled together imitating the composition (% relative abundance) of the original insect pools (Figure 1, “DNA pools”). In order to do this, an arbitrary 4 ng of DNA was used as the unit corresponding to 1 specimen and PCR products were diluted to a concentration of 4 ng/ $\mu$ L so to have a consistent volume for each pool (400  $\mu$ L) at a consistent concentration (4 ng/ $\mu$ L) (i.e., 4 ng of DNA for each insect composing the original pools). Libraries were prepared from these pools following the remainder of the metabarcoding protocol used for the whole insect pools.

To assemble mock communities representative of the post-amplification stage (Figure 1, “PCR Pools”), DNA extracted from each insect species as above was then amplified separately with each primer pair using the same PCR conditions as the whole insect mock communities. The PCR amplicons were then quantified using Qubit 2.0 Fluorometer (Figure 2) and diluted to a concentration of 2 ng/ $\mu$ L so to have a consistent volume for each pool (200  $\mu$ L) at a consistent

concentration (2 ng/μL), as outlined above (i.e., 2 ng of DNA for each insect composing the original pools, to reconstitute the original sample relative abundance), and libraries prepared as above. A second MiSeq run (“Run 2”) was performed three months after the first one (“Run 1”, which was conducted immediately) containing sequencing libraries generated from the DNA and PCR pools, as well as repeating the original insect pools using the same MiSeq machine and protocol. Run 2 of the insect pools was used to test how well the extracted DNA could be amplified and sequenced following an extended period of storage at -20°C, a typical temperature often used for long term storage of DNA samples, enabling future analysis on the same samples.

### Bioinformatics analysis

Raw sequence reads were demultiplexed using bcl2fastq v2.2.0 allowing for no mismatches to the expected index combinations (NCBI SRA acc no: PRJNA767112), then trimmed of PCR primer sequences using BB DuK v38.9 (Bushnell *et al.* 2017). Sequence quality profiles were used to further filter reads with >1 expected error (Edgar & Flyvbjerg 2015), or any ambiguous ‘N’ Bases, then remaining sequences were denoised using DADA2 v1.16 (Callahan *et al.* 2016) with the error model determined separately for each sequencing run. Following denoising, amplicon sequence variants (ASVs) inferred separately from each sequencing run were combined into a single table and any chimeric sequences removed de-novo using the “removeBimeraDenovo” function in DADA2. To further filter any non-specific amplification products and pseudogenes the ASV’s were aligned to a profile hidden Markov model (Eddy 1998) of the COI barcode region from Piper *et al.* (2021) using the aphid v1.3.3 R package (Wilkinson 2019), retaining sequences that met a minimum log-odds alignment score of 100 with a minimum match length of 100 bp. Retained ASVs were then checked for frame shifts and stop codons that commonly indicate pseudogenes (Roe & Sperling 2007). Taxonomy was determined by aligning ASVs to reference sequences of each taxon used in the mock communities using BLASTn v2.11.0 with a minimum percentage identity of 97% and minimum alignment coverage of 95%. The 52 % of ASVs (0.1% of abundance) that couldn’t be accurately mapped were discarded using filtering functions contained in the phyloseq v1.36.0 (McMurdie & Holmes 2013) and tidyverse v1.3.1 (Wickham *et al.* 2019) R packages.

### Statistical analysis

Differences in the number of ASVs detected between DNA extraction protocols, primer sets, and mock communities representing each workflow step were tested for significance using Analysis of Variance (ANOVA), followed by post-hoc pairwise comparisons with Tukey’s Honest Significant Difference (HSD) test (Tukey 1977). In order to compare the number of ASVs between each protocol without the confounding effect of differing read depths all samples were rarefied to 100,000 reads before ANOVAs were conducted. Differences in overall quantitative performance between each protocol was measured using the Root Mean Square Error (RMSE) between observed and expected relative abundances, an accuracy metric where smaller values indicate a less biased protocol. Taxonomic bias for each primer and workflow

step (Insect pools, DNA pools, amplicon pools) was estimated using a linear model of the compositional error (ratio between expected and observed abundances), with the results geometrically centered to be relative to the ‘average’ taxon (McLaren *et al.* 2018) and standard errors for each taxon coefficient generated from 1000 bootstrap resamples. The fit of the bias model for the dataset was evaluated by its ability to predict the observed relative abundances from sequencing using the known relative abundances in the mock communities. The separate bias estimates from each workflow step were then used to partition the total protocol bias (that represented by the Insect pools) into 3 different components for each primer set 1) DNA extraction, 2) PCR amplification, and 3) sequencing and bioinformatics as per McLaren *et al.* (2018). The change in abundance of each mock community taxon throughout the metabarcoding protocols relative to its starting abundance (Figure SM1) was then calculated by sequentially multiplying a starting abundance of 1 by the partitioned bias components for that taxon. Finally, the effects of morphological traits (biomass, sclerotization) and molecular traits (GC% of the whole amplicon, primer-template mismatch) on both the partitioned bias estimates and total protocol bias (represented by the Insect pools) was tested for significance using a second linear model fit to the bootstrapped bias estimates, and the relative influence of each trait on detection efficiency determined from its coefficient in the regression model. All statistical analyses were conducted within R 4.0.2 (R Core Team, 2020) using tidymodels v0.1.3 (Kuhn & Wickham 2020) packages.

## Results

### Comparison between non-destructive DNA extraction methods

A total of 7,241,574 and 8,676,552 reads were generated from the first and second MiSeq sequencing runs respectively, consisting of 182 unique ASVs. Each sample received a mean 272,600 ( $\pm 3,762$ ) reads, ranging from 615,687 reads for the highest sample to 619 reads for the lowest negative control sample, and contained a mean of 26.7 ASVs ( $\pm 0.119$ , range: 11-34). When reads that could not be classified to the known mock community members were removed (e.g., chimeras), the mean reads per sample dropped slightly to 263,581 ( $\pm 4,076$ ; range: 619-614,641) and the mean number of ASVs dropped to a mean 14.9 per sample ( $\pm 0.03$ , range: 7-16). Significant differences in the number of inferred ASVs were found between the different protocols (ANOVA;  $F_{(3, 60)} = 5.11$ ,  $p = .003$ ), primarily driven by the DNEasy treatments showing significantly more unique ASVs than the QuickExtract (Tukey’s HSD;  $p = .002$ ). On further exploration this difference was found to be due to a substantial dropout of taxa seen in the replicated QuickExtract samples on Run 2, which was run three months later (Figure SM1). When these analyses were repeated without the replicated QuickExtract samples, significant differences were again found between the treatments (ANOVA;  $F_{(3, 52)} = 9.78$ ,  $p < .001$ ), driven by both the QuickExtract and DNEasy treatments having higher numbers of ASVs than the DNA



pools (Tukey's HSD;  $p = .006$ ) and the PCR pools ( $p = .002$ ) with no significant difference between the QuickExtract and DNEasy samples ( $p > 0.05$ ). When considering only ASVs that could be classified to species level, significant differences remained between treatments (ANOVA;  $F(3, 52) = 6.18$ ,  $p = .001$ ), which pairwise comparisons revealed to be driven by the QuickExtract having significantly higher ASVs than the DNA (Tukey's HSD;  $p = .002$ ) or PCR pools ( $p = .001$ ), while no significant differences were found between the DNEasy and any other treatment ( $p > 0.05$ ). In contrast to the differences in protocols and pools, there were no significant differences between the two primers sets in terms of both total ASVs (ANOVA;  $F_{(1, 54)} = 0.21$ ,  $p > 0.05$ ) or ASVs that could be classified to species level (ANOVA;  $F_{(1, 54)} = 4.04$ ,  $p > 0.05$ ).

### Identification of the sources of bias

Both QuickExtract and DNEasy showed the highest deviation between expected and observed (from sequencing) relative abundances at 16% and 17% RMSE respectively for both primer pairs (Figure 3A) when compared to DNA pools (9% for the fwH2-fwH2n primer pair, and 7% for the fwH2-HexCOI4 primer pair) and the PCR pools (7% and 2%). There was a significant association found between the compositional error (ratio between expected and observed abundances) and species identity across all treatments and primer sets (Supplementary Table 1). When the taxon coefficients estimated by the bias model (Supplementary Table 2) were used to predict the observed relative abundances from the known number of individuals in the mock communities, a substantial improvement in quantitative performance was seen compared to the uncorrected sequencing data (Figure 3A). The model showed a good fit to the DNEasy pools, reducing the RMSE to 2% and 3% for the fwH2-fwH2n and fwH2-HexCOI4 primers respectively, and a near perfect fit to the DNA and PCR pools with 1% RMSE for both primers. While the bias model also reduced the RMSE for the QuickExtract samples, there was still substantial variance seen with an RMSE of 9% for fwH2-fwH2n, and 7% for the fwH2-HexCOI4 (Figure SM2).

### Taxon-specific biases

The taxonomic biases estimated by the model (Supplementary Table 2) revealed *A. solanicola*, *A. alternata*, *Acizzia* sp. and *R. padi* having the highest efficiency in the QuickExtract and DNEasy pools, and to a lesser extent the DNA pool. On the other hand, the two *Carpophilus* taxa showed the lowest relative efficiency (Figure 3B). The bias estimates for the QuickExtract insect pools showed substantially higher variance compared to the DNEasy insect pools across all taxa, with *Carpophilus truncatus* showing the largest difference between the two extraction methods (Figure 3B). When the total protocol bias was partitioned into the contribution of each step (Supplementary Table 3), marked differences in bias were seen between protocol steps on the same taxon (Figure 4A). For instance, *Carpophilus truncatus* saw substantially lower DNA extraction efficiency than the average taxon, but higher than average at the PCR stage. In contrast, *Diuraphis noxia* showed a higher DNA extraction efficiency at the DNA extraction

stage, but a lower efficiency in the later PCR step. The results for the PCR amplification are similarly reflected in the DNA concentrations obtained when single species were amplified to obtain mock communities representing post-PCR processes (Figure 2). Here we see that when individual PCRs are conducted from the same starting DNA concentration, *Diuraphis noxia* has much lower and *Carpophilus truncatus* has much higher DNA concentrations than the majority of the other taxa after both 30 and 40 cycles of amplification. The competing effects of different workflow stages on overall detection efficiency is particularly apparent when looking at the change in abundances of molecules for each taxon throughout the workflow (Figure 4B). Some taxa such as the psyllids *A. solanicola* and *A. alternata* consistently increase in abundance throughout the workflow, while others such as the fruit fly *B. tryoni* saw an initial increase in abundance, followed by a fall to almost equilibrium with the starting relative abundance (Figure 4B).

### Hardness- and biomass-associated bias

Differences in whole protocol bias was significantly associated with insect traits for both the QuickExtract ( $F_{(5,54)}=5.68$ ,  $R^2_{adj}=0.26$ ,  $p<0.001$ ) and DNEasy ( $F_{(5,58)}=20.50$ , adjusted  $R^2=0.60$ ,  $p<0.001$ ) protocols (Figure 5A). The efficiency with which a species was detected in a sample was positively influenced by whether that species had a soft (QuickExtract;  $\beta = 35.53$ , 95% CI [11.86, 22.10], DNEasy;  $\beta = 15.88$ , 95% CI [11.86, 22.10]), or intermediate hardness exoskeleton (QuickExtract;  $\beta = 19.13$ , 95% CI [6.09, 62.83], DNEasy;  $\beta = 8.55$ , 95% CI [6.65, 11.46]), or high amplicon GC% (QuickExtract;  $\beta = 31.93$ , 95% CI [7.53, 82.098], DNEasy;  $\beta = 13.71$ , 95% CI [9.90, 19.57])(Figure 5A). Insects with a hard exoskeleton showed no increase in detection efficiency for the DNEasy protocol ( $\beta = 0.12$ , 95% CI [-0.68, 0.86]), but increased efficiency for the QuickExtract protocol although with a confidence interval overlapping zero ( $\beta = 2.29$ , 95% CI [-0.43, 8.99]). Insect biomass unexpectedly had a strong negative effect on detection efficiency for both protocols (QuickExtract;  $\beta = -12.02$ , 95% CI [-41.50, -3.07], DNEasy;  $\beta = -4.73$ , 95% CI [-6.87, -3.40]), exceeded only by primer mismatch (QuickExtract;  $\beta = -23.20$ , 95% CI [-82.01, -5.84], DNEasy;  $\beta = -8.50$ , 95% CI [-12.66, -5.63])(Figure 5A). Generally, biases associated with each trait were more unpredictable (i.e., with larger variances observed) when using QuickExtract ( $R^2_{adj} = 0.26$ ) compared with DNEasy ( $R^2_{adj} = 0.60$ ) extraction method (Figure 5A). When considering just the bias contributed by the DNA extraction stage (Figure 5B), insect traits were also associated with detection efficiency albeit to a lesser degree than the whole protocol (QuickExtract;  $F_{(5,54)} = 8.79$ ,  $R^2_{adj} = 0.38$ ,  $p<0.001$ , DNEasy;  $F_{(5,58)} = 32.60$ ,  $R^2_{adj} = 0.70$ ,  $p<0.001$ ). Most notably, the relationship to insect biomass was reversed when considering just the DNA extraction stage, with larger biomass showing a slight positive effect on detection efficiency for both protocols (QuickExtract;  $\beta = 0.94$ , 95% CI [0.20, 1.81], DNEasy;  $\beta = 2.76$ , 95% CI [2.15, 3.51]) (Figure 5B). A soft (QuickExtract;  $\beta = 1.78$ , 95% CI [1.11, 2.97], DNEasy;  $\beta = 1.85$ , 95% CI [1.56, 2.20]) or intermediate hardness exoskeleton (QuickExtract;  $\beta = 2.56$ , 95% CI [1.85, 3.51], DNEasy;  $\beta = 2.79$ , 95% CI [2.50, 3.10]) again increased detection efficiency during DNA extraction, while a hard exoskeleton,

primer mismatch, or amplicon GC% all had confidence intervals overlapping zero for both protocols (Figure 5B). When considering only the bias contributed by the PCR stage, the overall effect of insect traits on detection efficiency was more comparable between protocols (QuickExtract  $F_{(5,54)}=7.13$ ,  $R^2_{adj} = 0.34$ ,  $p<0.001$ , DNEasy;  $F_{(5,58)}=7.65$ ,  $R^2_{adj} = 0.35$ ,  $p<0.001$ ) (Figure 5C). As expected, primer mismatch showed the strongest negative effect on detection efficiency (QuickExtract;  $\beta = -3.18$ , 95% CI [-3.08, -2.83], DNEasy;  $\beta = -3.08$ , 95% CI [-3.44, -2.83]), while amplicon GC% increased detection efficiency (QuickExtract;  $\beta = 3.15$ , 95% CI [2.65, 3.83], DNEasy;  $\beta = 3.14$ , 95% CI [2.86, 3.50]). All values of exoskeleton hardness showed a slight positive effect on detection efficiency at the PCR stage across both protocols ( $\beta = 2.37 - 3.09$ ), while biomass decreased efficiency (QuickExtract;  $\beta = -0.91$ , 95% CI [-1.16, -0.67], DNEasy;  $\beta = -0.92$ , 95% CI [-1.06, -0.79]).

## Discussion

Here, using non-destructive metabarcoding approaches, we successfully recorded all insect species present in pools composed of 100-101 individuals, including those species that were represented by just a single individual insect. At the same time, the use of these non-destructive DNA extraction methods allowed morphological voucher specimens of the insects to be preserved, as previously demonstrated (Martoni *et al.* 2019, Batovska *et al.* 2021). This is of paramount importance for many regulatory applications of metabarcoding, where retaining voucher specimens of potential pests or indicator taxa can be required for legal reasons, or to simply provide a morphological specimen preserved in an entomological collection for future taxonomic investigation (i.e., Martins *et al.* 2019, Martoni *et al.* 2019, Martoni *et al.* 2021a).

The results presented here show that non-destructive metabarcoding analysis can be successfully applied to bulk samples of agriculturally relevant insects to obtain a species identification, in agreement with recent studies (Carew *et al.* 2018; Nielsen *et al.* 2019; Batovska *et al.* 2021). This suggests that non-destructive metabarcoding has potential applications not only for biodiversity assessments but also for diagnostics and biosecurity purposes, to determine presence/absence of pests in bulk traps. Our study further highlights the molecular and morphological traits that affect quantification of individual insect species within non-destructively extracted bulk insect pools.

### Metabarcoding bias and non-destructive DNA extractions

The output generated by metabarcoding analyses is compositional data (Gloor *et al.* 2017). This means the relationship between the starting total abundance of a species and the output counts of sequence reads is completely lost, and the sequence reads returned for a taxon are only meaningful relative to the rest of the taxa within its sample (Gloor *et al.* 2017). Furthermore, the relationship between counts of sequence reads and the individuals they arise from is affected by a number of biases that systematically distort the measured sequence counts of each species from their true abundances (McLaren *et al.* 2019). This bias can lead to

taxonomic dropouts by diluting certain taxa below the detection limit and brings additional challenges for quantifying the number of individual specimens in a sample from metabarcoding sequences alone.

Amongst the known sources of metabarcoding bias, a range of physical characteristics of the insect community under study plays a key role, with perhaps the most obvious of these being the large variation in specimen body sizes within insect community assemblages (Chown & Gaston 2010). For example, even assuming that all available DNA is always extracted from a sample, the amount of DNA (and the reads obtained from it) from a single insect will depend on its biomass, which in turn depends on its species or life stage. Consequently, when aiming to use metabarcoding as a semi-quantitative technique, it is important to remember that the reads obtained from a single insect will vary depending on a number of its characteristics (Thomas *et al.* 2016; McLaren *et al.* 2019), as well as on the composition of the insect pool analysed (Gloor *et al.* 2017).

With non-destructive DNA extractions, an additional factor to consider is that non-destructive protocols mostly act upon the external surface of the insect exoskeletons, as opposed to destructive methods that can potentially access all of the DNA contained in the insect bodies. Therefore, when comparing the number of reads obtained from the insect pools with that from the DNA pools two additional factors had to be taken into consideration in addition to sample composition: exoskeleton robustness and species biomass. When comparing these factors, a soft exoskeleton had the strongest positive effect on relative efficiency, suggesting that soft-bodied insects facilitate non-destructive DNA extraction, resulting in a higher representation than one would expect from their relative biomass within the sample. Whether this proves beneficial or problematic for species detection will depend on the specific communities a surveillance programme is targeting. On the one hand, the smaller influence of species biomass could alleviate the requirement for any morphological pre-sorting to ensure large specimens don't drown out smaller ones (Elbrecht *et al.* 2017). However, the negative effect of hard exoskeletons on detection efficiency could produce false negatives for low abundance taxa with a high level of sclerotization. This is the most likely explanation for the differences in extraction efficiency between taxa such as *Carpophilus truncatus*, having a highly sclerotized exoskeleton, in comparison to the more soft-bodied insects such as *Diuraphis noxia* (Figure 5A).

These results were consistent between both the QuickExtract and DNEasy kit protocols, suggesting that the bias towards softer bodied insects is an inherent aspect of the non-destructive DNA extraction process itself, rather than a specific protocol or kit. Nevertheless, there are a few key differences between the two extraction methods we evaluated that will influence their practical application within future studies. Firstly, the QuickExtract method was substantially faster, requiring only one hour of operator's time as opposed to the overnight incubation of the DNEasy method. However, QuickExtract produced unpredictable variation in DNA extraction performance and appeared less suitable for long-term DNA preservation. When the DNA extracts from the insect mock communities were re-sequenced three months after extraction – during which the templates were kept in -20 °C freezer – the QuickExtract template appeared to

have degraded dramatically, resulting in a number of species dropping out. On the other hand, the presence/absence of species from the DNEasy kit was virtually identical to when they were sequenced within days of being extracted. Therefore, DNA extraction products obtained with the QuickExtract kit should be used immediately for analysis, or stored at lower temperatures, i.e., in a -80 °C freezer (as suggested in the manufacturer's recommendations) if possible.

In addition, the results showed relative abundance measurements obtained through non-destructive metabarcoding assays are most strongly influenced by the DNA extraction processes. DNA extraction was the not just the largest contributor to protocol bias, but also had the highest variance of all the tested workflow stages, making it less predictable. Nevertheless, it is unclear if this aspect is specific to the non-destructive DNA extraction process or is applicable to DNA extractions generally. Few metabarcoding studies have attempted to partition the total protocol bias into separate components for each major protocol step. It would be valuable for future research to compare the quantitative results of destructive and non-destructive extractions, using a similar bias-partitioning approach to that implemented within our study. The non-destructive DNA extraction bias issues highlighted here may turn out to be unavoidable when voucher specimens are required, as destructive DNA extraction is not an option.

### **Primer mismatch and PCR bias**

Beside DNA extraction, we wanted to assess if bias was introduced also at the PCR stage. Mismatches between PCR primers and template molecules have been considered the primary contributor to metabarcoding bias, particularly at the 3' end of the primer where nucleic acid extension takes place (Piñol *et al.* 2015). Primer-template mismatch can be particularly problematic for protein coding genes such as COI, due to natural degeneracy in the genetic code leaving no strictly conserved gene regions for design of universal PCR primers (Deagle *et al.* 2014). This has necessitated the inclusion of multiple degenerate nucleotide bases in metabarcoding primers in order to mitigate the effects of mismatch on detection efficiency and abundance estimates (Elbrecht *et al.* 2019).

While much of the previous literature has highlighted that bias introduced during PCR amplification is responsible for the semi-quantitative nature of metabarcoding results, our study observed a smaller contribution of PCR compared to the DNA extraction process. While this could potentially be attributed to non-destructive DNA extraction being more bias prone than its homogenisation-based alternatives, this will require further experiments to demonstrate. Alternatively, the smaller contribution of the PCR process to total protocol bias could have been due to the two primer pairs employed in this study containing a number of degenerate bases, having been designed to be generic for insects (Vamos *et al.* 2017; Marquina *et al.* 2018; Elbrecht *et al.* 2019). Previous studies that have found high levels of PCR bias have mostly used less-degenerate primers or primers which have since been identified to contain critical design flaws (Elbrecht & Leese 2017; Piper *et al.* 2021). Nevertheless, some species included in our study still showed mismatch to these generic primers (Table 1), and this was a primary driver of reduced efficiency at the PCR stage (Figure 4C).

Additionally, primer bias was confirmed when quantifying the PCR product obtained from each insect species in separate PCRs, where each primer pair amplified the DNA of some species up to ten times more than others. Since this DNA had been normalized prior to PCR, variation in subsequent concentration could be linked bias introduced during amplification. Such variation in the concentration appeared to also depend on the number of PCR cycles performed, with a higher number of cycles showing more similar concentrations across different insect species. Nevertheless, this bias did not impact the presence/absence of targets, with both primer pairs tested here recording all of the species present in all the pools, with a sensitivity of up to 1 in 101 for many of the species tested. Further work could be required to test if primer bias may have greater effects on presence/absence in larger communities, or at a lower sequencing depth, as well as if additional primer pairs might produce a lower primer bias. Nonetheless, it is important to remember that, in a diagnostic context, priority is given to a precise presence/absence assessment, especially when testing for the presence of unexpected pests or rare species. Therefore, for general insect biodiversity monitoring and surveillance, the semi-quantitative but taxonomically broader screening of agricultural traps or biodiversity assessment samples that can be performed using MiSeq metabarcoding, could be considered the preferred approach over a more quantitative but perhaps more time-consuming method, such as qPCR, targeting each individual species separately.

#### **Possible strategies to account for metabarcoding bias.**

When considering the number of reads generated as output of the MiSeq metabarcoding, the relative abundance of reads obtained for each insect species was correlated with the number of individuals present in the pool but biased toward certain species. Following the multiplicative model of McLaren *et al.* (2019), we used mock communities representing each step of the library preparation pipeline to partition the total protocol bias into that explained by each component. Indeed, the relative abundances obtained from different insect species appeared to be subject to bias during the different steps of the laboratory workflow, including DNA extraction and PCR, fundamentally due to both morphological and molecular traits differing across each insect group. Amongst these, some of the morphological traits determined here to significantly affect the metabarcoding results are surface, size and consistency of the insect tissues, as previously reported for Coleoptera (Martins *et al.* 2019). On the other hand, molecular traits include PCR primers mismatch, which was observed here for both the primer pairs used, despite these being considered generic primers (Vamos *et al.* 2017; Marquina *et al.* 2018). In the literature, attempts have been made to improve a quantitative output for metabarcoding, either by developing new primers, or protocols (e.g., Elbrecht & Leese 2015; Deagle *et al.* 2019; Lamb *et al.* 2019). Our results indicate that understanding the individual contributions of each laboratory stage, rather than just the total overall protocol bias, should be a critical consideration for future efforts. This was particularly apparent when measuring the effect of biomass on detection efficiency, as when the total protocol bias was considered, larger biomass appeared to have a strong negative effect on the number of reads for each given species. On the other hand, when considering just the bias

introduced by the DNA extraction stage, larger insect biomass increased the number of reads produced, as expected from previous studies. Therefore, when optimizing protocols researchers should be wary of the often-contrasting effects of the different laboratory steps, in order to not be confounded into mistakenly attempting to optimize one aspect, when in reality the majority of the bias is being introduced through a different laboratory step.

While efforts to further optimise protocols will no doubt prove important for increasing the quantitative performance of metabarcoding, a less explored but complementary approach is to use statistical models to actively correct for bias during analysis (Thomas *et al.* 2016; Krehenwinkel *et al.* 2017; McLaren *et al.* 2019). Promisingly, we found bias to be highly predictable with our model, demonstrating the potential for developing correction factors to improve the quantitative results of metabarcoding assay for specific important target species, e.g., agricultural pests. However, even with our consistently treated mock communities, DNA extraction bias had the highest variance, and thus would be the most difficult factor to model and correct for. In a surveillance situation, this could be exacerbated by differences in environmental conditions between real trap samples and the mock communities used to derive the correction factors introducing further bias (Krehenwinkel *et al.* 2018). Furthermore, developing correction factors from mock communities requires knowing the organisms that will likely be encountered *a priori*, as well as being able to acquire specimens of them. Unfortunately, in both contexts of biodiversity monitoring and surveillance, knowing *a priori* what organisms will be recorded in order to prepare targeted spike-ins or tailored mock community could be challenging and, often, defy the very purpose of the study.

## Conclusions

Based on the results obtained here, different insect species, even belonging to the same order or genus, are subject to different degrees of bias during a non-destructive metabarcoding workflow. These biases are driven by species-specific morphological and molecular traits that interact with different protocol steps to either increase or decrease detection efficiency. If bias can be measured for that taxon *a priori* using mock communities, correction factors could therefore be used to calibrate the results to better reflect the actual abundances, but requires knowledge of the sample composition prior to analysis.

Ultimately, if the surveillance program or the biodiversity assessment are conducted long-term using a non-destructive DNA extraction and thus retaining voucher specimens, obtaining specimens to develop correction factors covering a wider diversity may eventually be achievable. Inclusion of mock communities in the metabarcoding analysis would not only be possible, but would sensibly improve in accuracy based on an ongoing monitoring of the geographical area of interest. This highlights the importance of long-term monitoring and biodiversity data collection projects not only for agricultural and industrial areas but also for adjacent ecosystems that can contribute to the species diversity recorded during insect trapping programs.

## Acknowledgements

The authors thank PeerJ editor, Dr Sheila Colla, and two anonymous reviewers for their useful comments and suggestions. The authors would like to thank the iMapPESTS project, supported by Horticulture Innovation Australia (ST16010) through funding from the Australian Government Department of Agriculture as part of its Rural R&D for Profit program and Grains Research and Development Corporation, for funding these experiments and the publication of this article. We thank Isabel Valenzuela, Jessi Henneken, and Will Boston (Agriculture Victoria) for providing many of the laboratory colony insect specimens used in this study.

## Data availability

Raw sequence data is available at NCBI SRA acc no: PRJNA767112. All code and additional files required to reproduce the bioinformatic and statistical analyses presented in this manuscript are available at the following GitHub repository:  
[https://github.com/alexpiper/disentangling\\_metabarcoding\\_bias\\_MS](https://github.com/alexpiper/disentangling_metabarcoding_bias_MS)



# References

- Andújar C, Arribas P, Yu DW, Vogler AP & Emerson BC (2018) Why the COI barcode should be the community DNA metabarcode for the metazoa. *Molecular Ecology*, 27: 3968–3975. <https://doi.org/10.1111/mec.14844>.
- Armstrong KF. & Ball SL. (2005) DNA barcodes for biosecurity: invasive species identification. *Philosophical transactions of the Royal Society B*, 360(1462): 1813–1823. <https://doi.org/10.1098/rstb.2005.1713>.
- Ashfaq M, Hebert PDN & Naum A. (2016) DNA barcodes for bio-surveillance: Regulated and economically important arthropod plant pests. *Genome*, 59(11): 933–945. <https://doi.org/10.1139/gen-2016-0024>.
- Bahder BW, Bollinger L, Sudarshana MR & Zalom FG. (2015) Preparation of mealybugs (Hemiptera: Pseudococcidae) for genetic characterization and morphological examination. *Journal of Insect Science*, 15: 104. <https://doi.org/10.1093/jisesa/iev086>.
- Batovska J, Piper AM, Valenzuela I, Cunningham JP & Blacket MJ. (2021) Developing a Non-destructive Metabarcoding Protocol for Detection of Pest Insects in Bulk Trap Catches. *Scientific Reports*, 11, 7946. doi:10.1038/s41598-021-85855-6
- Braukmann TW, Ivanova NV, Prosser SW, Elbrecht V, Steinke D, Ratnasingham S, de Waard JR, Sones JE, Zakharov EV & Hebert PDN (2019) Metabarcoding a diverse arthropod mock community. *Molecular ecology resources*, 19(3): 711–727. <https://doi.org/10.1111/1755-0998.13008>.
- Bushnell B, Rood J, & Singer E. (2017) BBMerge – Accurate paired shotgun read merging via overlap. *PLOS ONE*, 12: 1–15. <https://doi.org/10.1371/journal.pone.0185056>.
- Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP (2016) DADA2: High-resolution sample inference from Illumina amplicon data. *Nature Methods*, 13: 581–583. <https://doi.org/10.1038/nmeth.3869>.
- Carew ME, Coleman RA & Hoffmann AA (2018) Can non-destructive DNA extraction of bulk invertebrate samples be used for metabarcoding? *PeerJ*, 6: e4980. <https://doi.org/10.7717/peerj.4980>.
- Chown SL, Hodgins KA, Griffin PC, Oakeshott, J.G., Byrne, M. & Hoffmann, A.A. (2014) Biological invasions, climate change and genomics. *Evolutionary Applications*, 8(1): 23–46. <https://doi.org/10.1111/eva.12234>.
- Chown SL & Gaston KJ. (2010) Body size variation in insects: a macroecological perspective. *Biological Reviews*, 85(1): 139–169. <https://doi.org/10.1111/j.1469-185X.2009.00097.x>.
- Darling JA, Pochon X, Abbott CL, Inglis GJ & Zaiko A. (2020) The risks of using molecular biodiversity data for incidental detection of species of concern. *Diversity and Distribution*, 26(9): 1116–1121. <https://doi.org/10.1111/ddi.13108>.

- Davidovitch L, Stoklosa R, Majer J, Nietrzeba A, Whittle P, Mengersen K & Ben-Haim Y. (2009) Info-gap theory and robust design of surveillance for invasive species: The case study of Barrow Island. *Journal of Environmental Management*, 90(8): 2785–2793. <https://doi.org/10.1016/j.jenvman.2009.03.011>.
- Deagle BE., Jarman SN, Coissac E, Pompanon F & Taberlet P. (2014) DNA metabarcoding and the cytochrome c oxidase subunit I marker: Not a perfect match. *Biology Letters*, 10(9), 20140562. <https://doi.org/10.1098/rsbl.2014.0562>.
- Eddy SR. (1998) Profile hidden Markov models. *Bioinformatics* (Oxford, England), 14(9): 755–763.
- Edgar RC & Flyvbjerg H. (2015) Error filtering, pair assembly and error correction for next-generation sequencing reads. *Bioinformatics*, 31(21): 3476–3482. <https://doi.org/10.1093/bioinformatics/btv4>.
- Elbrecht, V., & Leese, F. (2017). Validation and Development of COI Metabarcoding Primers for Freshwater Macroinvertebrate Bioassessment. *Frontiers in Environmental Science*, 5, 11. <https://doi.org/10.3389/fenvs.2017.00011>
- 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). <https://doi.org/10.1371/journal.pone.0130324>.
- Elbrecht V, Peinert B & Leese F. (2017). Sorting things out: Assessing effects of unequal specimen biomass on DNA metabarcoding. *Ecology and evolution*, 7(17): 6918–6926. <https://doi.org/10.1002/ece3.3192>.
- Elbrecht V, Braukmann TWA, Ivanova NV, Prosser SWJ, Hajibabaei M, Wright M, Zakharov EV, Hebert PDN & Steinke D. (2019) Validation of COI metabarcoding primers for terrestrial arthropods. *PeerJ*, 7, e7745. <https://doi.org/10.7717/peerj.7745>.
- EPPO. (2021). PM 7/129 (2) DNA barcoding as an identification tool for a number of regulated pests. *EPPO Bulletin*, 51(1), 100–143. <https://doi.org/10.1111/epp.12724>
- Girón JC & Short AEZ (2021) The Acidocerinae (Coleoptera, Hydrophilidae): taxonomy, classification, and catalog of species. *ZooKeys*, 1045: 1–236. <https://doi.org/10.3897/zookeys.1045.63810>.
- Gloor GB, Macklaim JM, Pawlowsky-Glahn V & Egozcú JJ (2017) Microbiome datasets are compositional: and this is not optional. *Frontiers in microbiology*, 8: 2224. <https://doi.org/10.3389/fmicb.2017.02224>.
- Hebert PDN, Cywinska A, Ball SL & DeWaard JR (2003) Biological identifications through DNA barcodes. *Proceedings of the Royal Society B - Biological Sciences*, 270: 313–321. <https://doi.org/10.1098/rspb.2002.2218>.
- Hulme PE (2009) Trade, transport and trouble: Managing invasive species pathways in an era of globalization. *Journal of Applied Ecology*, 46: 10–18. <https://doi.org/10.1111/j.1365-2664.2008.01600.x>.

- 649 • Krehenwinkel H, Wolf M, Lim JY, Rominger AJ, Simison WB & Gillespie RG. (2017)

650 Estimating and mitigating amplification bias in qualitative and quantitative arthropod

651 metabarcoding. *Scientific Reports*, 7, 17668. [https://doi.org/10.1038/s41598-017-17333-](https://doi.org/10.1038/s41598-017-17333-x)

652 [x](https://doi.org/10.1038/s41598-017-17333-x).
- 653 • Krehenwinkel H, Fong M, Kennedy S, Huang EG, Noriyuki S, Cayetano L & Gillespie

654 R. (2018) The effect of DNA degradation bias in passive sampling devices on

655 metabarcoding studies of arthropod communities and their associated microbiota. *PLOS*

656 *ONE*, 13(1), e0189188. <https://doi.org/10.1371/journal.pone.0189188>.
- 657 • Kuhn M & Wickham H. (2020) Tidymodels: a collection of packages for modeling and

658 machine learning using tidyverse principles. Available at: <https://www.tidymodels.org>.
- 659 • Lamb PD, Hunter E, Pinnegar JK, Creer S, Davies RG & Taylor MI. (2019) How

660 quantitative is metabarcoding: A meta-analytical approach. *Molecular ecology*, 28(2):

661 420–430. <https://doi.org/10.1111/mec.14920>.
- 662 • Lefort M-C, Beggs JR, Glare TR, Saunders TE, Doyle EJ & Boyer S (2020) A molecular

663 approach to study Hymenoptera diets using wasp nests. *NeoBiota*, 63: 57–79.

664 <https://doi.org/10.3897/neobiota.63.58640>.
- 665 • Low-Choy S. (2015) *Getting the story straight: Laying the foundations for statistical*

666 *evaluation of the performance of surveillance*. In: Jarrad F, Low-Choy S, Mengersen K,

667 eds. *Biosecurity Surveillance: Quantitative Approaches*. 6th ed. CABI; 43–73.
- 668 • Marquina D, Andersson AF & Ronquist, F. (2018) New mitochondrial primers for

669 metabarcoding of insects, designed and evaluated using in silico methods. *Molecular*

670 *Ecology Resources*, 19(1): 90–104. <https://doi.org/10.1111/1755-0998.12942>.
- 671 • Marquina D, Esparza-Salas R, Roslin T & Ronquist F. (2019) Establishing arthropod

672 community composition using metabarcoding: Surprising inconsistencies between soil

673 samples and preservative ethanol and homogenate from Malaise trap catches. *Molecular*

674 *Ecology Resources*, 19(6): 1516–1530. <https://doi.org/10.1111/1755-0998.13071>.
- 675 • Martins FMS, Galhardo M, Filipe AF, Teixeira A, Pinheiro P, Paupério J, Alves PC &

676 Beja P. (2019) Have the cake and eat it: Optimizing nondestructive DNA metabarcoding

677 of macroinvertebrate samples for freshwater biomonitoring. *Molecular ecology*

678 *resources*, 19(4): 863–876. <https://doi.org/10.1111/1755-0998.13012>.
- 679 • Martoni F, Valenzuela I & Blacket M. (2019) Non-destructive DNA extractions from fly

680 larvae (Diptera: Muscidae) enable molecular identification of species and enhance

681 morphological features. *Austral Entomology*, 56(4): 848–856.

682 <https://doi.org/10.1111/aen.12419>.
- 683 • Martoni F, Taylor GS, Blacket M. (2020) Illuminating insights into the biodiversity of the

684 Australian psyllids (Hemiptera: Psylloidea) collected using light trapping. *Insects*, 11(6):

685 354. <https://doi.org/10.3390/insects11060354>.
- 686 • Martoni F, Valenzuela I & Blacket M. (2021a) On the complementarity of DNA

687 barcoding and morphology to distinguish benign endemic insects from possible pests: the

case of *Dirioxa pornia* and the tribe Acanthonevrini (Diptera: Tephritidae: Phytalmiinae) in Australia. *Insect Science*, 28(1): 261-270 <https://doi.org/10.1111/1744-7917.12769>.

- Martoni F, Nogarotto E, Piper AM, Mann R, Valenzuela I, Eow L, Rako L, Rodoni BC, Blacket MJ. (2021b) Propylene Glycol and Non-Destructive DNA Extractions Enable Preservation and Isolation of Insect and Hosted Bacterial DNA. *Agriculture*. 2021; 11(1):77. <https://doi.org/10.3390/agriculture11010077>
- McLaren MR, Willis AD & Callahan BJ (2019) Consistent and correctable bias in metagenomic sequencing experiments. *eLife*, 8:e46923. <https://doi.org/10.7554/eLife.46923>.
- McMurdie PJ & Holmes S. (2013) phyloseq: An R Package for Reproducible Interactive Analysis and Graphics of Microbiome Census Data. *PLOS ONE*, 8, e61217. <https://doi.org/10.1371/journal.pone.0061217>.
- Meyerson LA & Mooney HA. (2007) Invasive alien species in an era of globalization. *Frontiers in Ecology and the Environment*, 5: 199–208. [https://doi.org/10.1890/1540-9295\(2007\)5\[199:IASIAE\]2.0.CO;2](https://doi.org/10.1890/1540-9295(2007)5[199:IASIAE]2.0.CO;2).
- Nielsen M, Gilbert MTP, Pape T, Bohmann K (2019) A simplified DNA extraction protocol for unsorted bulk arthropod samples that maintains exoskeletal integrity. *Environmental DNA*, 1(2), 144–154. <https://doi.org/10.1002/edn3.16>.
- Paini DR, Sheppard AW, Cook DC, De Barro, PJ, Worner S & Thomas MB (2016) Global threat to agriculture from invasive species. *Proceedings of the National Academy of Science USA*, 113: 7575–9. <https://doi.org/10.1073/pnas.1602205113>.
- 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(4): 819–830. <https://doi.org/10.1111/1755-0998.12355>.
- Piper A, Batovska J, Cogan NOI, Weiss J, Cunningham JP, Rodoni BC & Blacket MJ. (2019) Prospects and challenges of implementing DNA metabarcoding for high-throughput insect surveillance. *GigaScience*, 8: 1–22. <https://doi.org/10.1093/gigascience/giz092>.
- Piper AM, Cogan NOI, Cunningham JP & Blacket MJ. (2021) Computational Evaluation of DNA Metabarcoding for Universal Diagnostics of Invasive Insect Pests. *BioRxiv*. <https://doi.org/10.1101/2021.03.16.435710>.
- Pohjoismäki JLO, Karhunen PJ, Goebeler S, Saukko P & Sääksjärvi IE. (2010) Indoors forensic entomology: colonization of human remains in closed environments by specific species of sarcosaprophagous flies. *Forensic Science International*, 199 (1–3): 38–42. <https://doi.org/10.1016/j.forsciint.2010.02.033>.
- R Core Team (2020) R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

- 726 • Roe AD & Sperling FA. (2007) Patterns of evolution of mitochondrial cytochrome c  
727 oxidase I and II DNA and implications for DNA barcoding. *Molecular Phylogenetics and*  
728 *evolution*, 44(1): 325–345. <https://doi.org/10.1016/j.ympev.2006.12.005>.
- 729 • Schutze MK, Bourtzis K, Cameron SL, Clarke AR, De Meyer M, Hee AKW, Hendrichs  
730 J, Krosch MN & Mwatawala M. (2017). Integrative taxonomy versus taxonomic  
731 authority without peer review: the case of the Oriental fruit fly, *Bactrocera dorsalis*  
732 (Tephritidae). *Systematic Entomology*, 42: 609–620. <https://doi.org/10.1111/syen.12250>.
- 733 • Thomas AC, Deagle BE, Eveson JP, Harsch CH & Trites AW. (2016) Quantitative DNA  
734 metabarcoding: Improved estimates of species proportional biomass using correction  
735 factors derived from control material. *Molecular Ecology Resources*, 16(3): 714–726.  
736 <https://doi.org/10.1111/1755-0998.12490>.
- 737 • Tukey JW (1977) Exploratory data analysis. Addison-Wesley, Reading.
- 738 • Vamos EE, Elbrecht V & Leese F. (2017) Short COI markers for freshwater  
739 macroinvertebrate metabarcoding. *Metabarcoding and Metagenomics*, 1: e14625.  
740 <https://doi.org/10.3897/mbmg.1.14625>.
- 741 • Wickham H, Averick M, Bryan J, Chang W, D'Agostino McGowan L, François R,  
742 Grolemond G, Hayes A, Henry L, Hester J, Kuhn M, Pedersen TL, Miller E, Bache SM,  
743 Müller K, Ooms J, Robinson D, Seide DP, Spinu V, Takahashi K, Vaughan D, Wilke C,  
744 Woo K & Yutani, H.(2019). Welcome to the Tidyverse. *Journal of Open Source*  
745 *Software*, 4(43): 1686. <https://doi.org/10.21105/joss.01686>.
- 746 • Wilkinson, S. (2019). aphid: an R package for analysis with profile hidden Markov  
747 models. *Bioinformatics*, 35(19), 3829–3830.  
748 <https://doi.org/10.1093/bioinformatics/btz159>
- 749 • Yu DW, Ji Y, Emerson BC, Wang X, Ye C, Yang C & Ding Z. (2012) Biodiversity soup:  
750 metabarcoding of arthropods for rapid biodiversity assessment and biomonitoring.  
751 *Methods in Ecology and Evolution*, 3: 613–623. [https://doi.org/10.1111/j.2041-](https://doi.org/10.1111/j.2041-210X.2012.00198.x)  
752 [210X.2012.00198.x](https://doi.org/10.1111/j.2041-210X.2012.00198.x).

# **Figures**

**Figure 1:** Workflow of the experiment for the three types of pools. Insect pools are the result of non-destructive DNA extractions from the pooled insect specimens. DNA pools are the result of pooled DNA that was destructively extracted from each insect species. PCR pools are the result of DNA that was extracted, quantified, and amplified separately from each insect species then pooled before indexing qPCR. All pools were indexed, sequenced and analysed following the same protocol. Each pool was amplified using two primer sets.

**Figure 2:** DNA concentration measured from the PCR amplification obtained for each of the two primer pairs after 30 PCR cycles **(A)** or 40 PCR cycles **(B)** on 4 ng of DNA extracted from each single insect species.

**Figure 3: A)** Comparison of observed relative abundances to expected relative abundances of each taxon **B)** Estimated bias for each taxon displayed relative to the geometric mean bias, with 95% and 50% confidence intervals displayed. RMSE; Root Mean Square Error.

**Figure 4: A)** Metabarcoding bias partitioned to each protocol and primer step, with 95% and 50% confidence intervals displayed. Change in composition of each taxon throughout the workflow relative to the geometric mean taxon with standard errors displayed, for **B)** QuickExtract, and **C)** DNEasy obtained by sequentially multiplying a starting abundance of 1 by the partitioned bias estimates from panel A. Insect species represented in graphs B and C are arranged and colour-coded based on graph A.

**Figure 5: A)** Coefficients of model predicting the estimated taxon-specific bias from the traits in panel A for the unpartitioned bias and the 3 partitioned bias steps. Coefficients are displayed on a pseudo-log scale to avoid compressing around zero.

## **Supplementary figures**



**Figure SM1:** Comparison of quantitative **A)** and qualitative **B)** results for DNA extraction methods on whole insect pools in Run 1, sequenced immediately, compared with Run 2, after three-months storage at -20 °C.

**Figure SM2:** Fit of the bias model to each set of pools. Proportions are displayed on a pseudo-log scale to avoid compressing variation around zero. RMSE; Root Mean Square Error.

## Tables

**Table 1: Composition of the eight pools used for this study.** The number of individual insects is reported for each pool, as well as the total number of individuals (in bold). DNA and PCR pools were assembled with the same proportions reported here.

**Table 2: Molecular and morphological characteristics of the insects used in the pools.** Molecular traits reported include primer mismatch and GC% of the amplified sequence for each of the two primer pairs used. Morphological traits include measurements of insect volume, obtained by averaging measures across 10 specimens per species, as well as a scale of exoskeleton hardness. Hardness values are 1=soft, 2=intermediate, 3=hard.

## Supplementary Tables

**Supplementary Table 1: Overall model fits across the primers and mock community types.**

**Supplementary Table 2: Estimated coefficients (bias) for each taxon, with bootstrap standard errors across the primers and mock community types.**

**Supplementary Table 3: Partitioned bias, and bootstrap standard errors across the primers and mock community types.**

# **Table 1**(on next page)

Table 1: Composition of the eight pools used for this study.

The number of individual insects is reported for each pool, as well as the total number of individuals (in bold). DNA and PCR pools were assembled with the same proportions reported here.



- 1 **Table 1: Composition of the eight pools used for this study.** The number of
- 2 individual insects is reported for each pool, as well as the total number of individuals (in
- 3 bold). DNA and PCR pools were assembled with the same proportions reported here.

		Pool							
Species	Order	1	2	3	4	5	6	7	8
<i>Carpophilus davidsoni</i>	Coleoptera	25	50	5	10	25	50	5	10
<i>C. truncatus</i>	Coleoptera	1	1	1	1	1	1	1	1
<i>Bactrocera tryoni</i>	Diptera	1	1	1	1	1	1	1	1
<i>Bradysia</i> nr. <i>ocellaris</i>	Diptera	0	0	0	0	1	0	1	0
<i>Drosophila hydei</i>	Diptera	1	1	1	1	1	1	1	1
<i>D. melanogaster</i>	Diptera	5	10	25	50	5	10	25	50
<i>D. simulans</i>	Diptera	1	1	1	1	1	1	1	1
<i>Scaptodrosophila lativittata</i>	Diptera	1	1	1	1	1	1	1	1
<i>Acizzia alternata</i>	Hemiptera	1	1	1	1	1	1	1	1
<i>A. solanicola</i>	Hemiptera	10	25	50	5	10	25	50	5
<i>Acizzia</i> sp.	Hemiptera	1	1	1	1	1	1	1	1
<i>Diuraphis noxia</i>	Hemiptera	1	1	1	1	1	1	1	1
<i>Metopolophium dirhodum</i>	Hemiptera	1	1	1	1	1	1	1	1
<i>Rhopalosiphum padi</i>	Hemiptera	1	1	1	1	1	1	1	1
<i>Aphidius colemani</i>	Hymenoptera	36	5	6	12	36	4	9	20
<i>Lysiphlebus testaceipes</i>	Hymenoptera	14	0	4	13	14	1	1	5
Total individuals		100	100	100	100	101	100	101	100

4

5

## Table 2 (on next page)

Table 2: Molecular and morphological characteristics of the insects used in the pools.

Molecular traits reported include primer mismatch and GC% of the amplified sequence for each of the two primer pairs used. Morphological traits include measurements of insect volume, obtained by averaging measures across 10 specimens per species, as well as a scale of exoskeleton hardness. Hardness values are 1=soft, 2=intermediate, 3=hard.

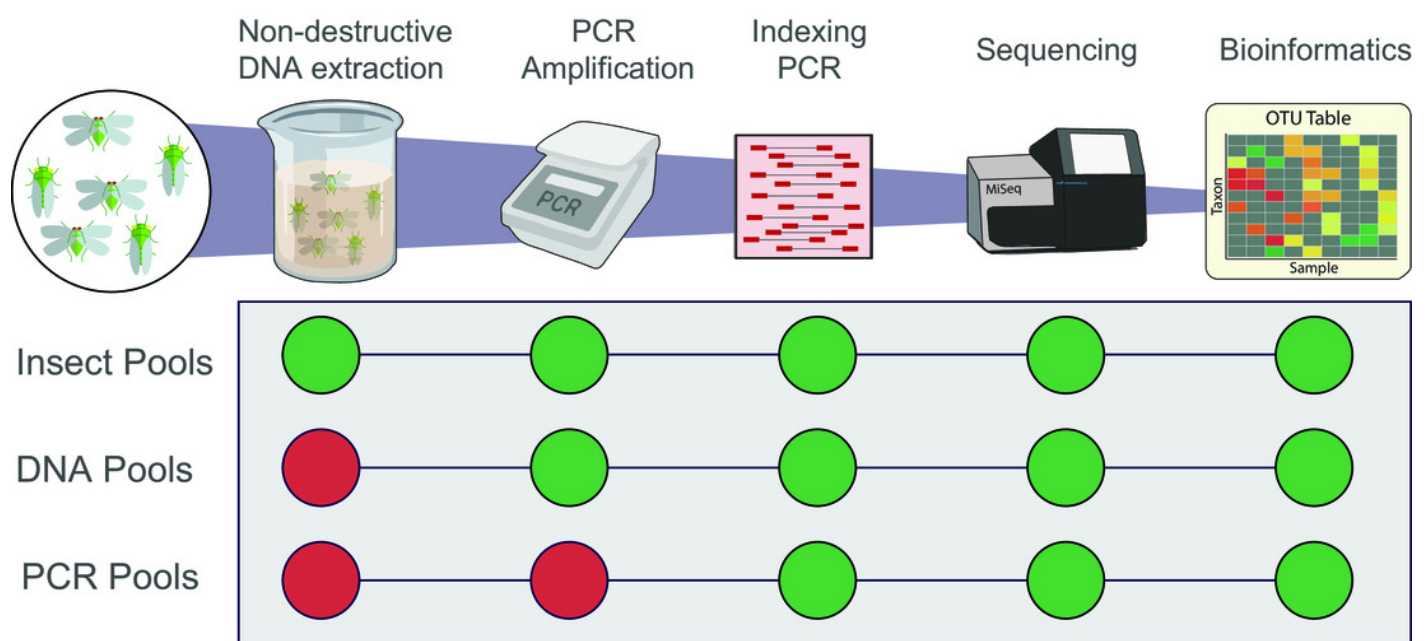
**Table 2: Molecular and morphological characteristics of the insects used in the pools.** Molecular traits reported include primer mismatch and GC% of the amplified sequence for each of the two primer pairs used. Morphological traits include measurements of insect volume, obtained by averaging measures across 10 specimens per species, as well as a scale of exoskeleton hardness. Hardness values are 1=soft, 2=intermediate, 3=hard.

Species	Order	Molecular traits						Morphological traits	
		fwhF2 mismatch	fwhR2n mismatch	Amplicon GC%	fwhF2 mismatch	HexCOIR4 mismatch	Amplicon GC%	Volume (mm <sup>3</sup> )	Hardness
<i>Carpophilus davidsoni</i>	Coleoptera	0.12	0.09	0.37	0.12	0.06	0.38	3.8	3
<i>Carpophilus truncatus</i>	Coleoptera	0.04	0.04	0.35	0.04	0	0.35	3.91	3
<i>Bactrocera tryoni</i>	Diptera	0	0.04	0.37	0	0	0.37	26.51	2
<i>Bradysia nr. ocellaris</i>	Diptera	0.08	0	0.34	0.08	0	0.34	0.25	2
<i>Drosophila hydei</i>	Diptera	0.04	0	0.33	0.04	0.06	0.33	2.19	2
<i>Drosophila melanogaster</i>	Diptera	0	0	0.3	0	0	0.3	1.14	2
<i>Drosophila simulans</i>	Diptera	0	0	0.32	0	0	0.32	0.79	2
<i>Scaptodrosophila lativittata</i>	Diptera	0	0	0.28	0	0	0.28	3.09	2
<i>Acizzia alternata</i>	Hemiptera	0	0.04	0.3	0	0	0.31	0.22	1
<i>Acizzia solanicola</i>	Hemiptera	0	0.04	0.32	0	0	0.33	0.47	1
<i>Acizzia sp.</i>	Hemiptera	0.08	0.09	0.32	0.08	0	0.33	1.18	1
<i>Diuraphis noxia</i>	Hemiptera	0.04	0	0.2	0.04	0	0.2	0.5	1
<i>Metopolophium dirhodum</i>	Hemiptera	0.04	0	0.22	0.04	0	0.22	0.65	1
<i>Rhopalosiphum padi</i>	Hemiptera	0.04	0	0.19	0.04	0	0.2	0.18	1
<i>Aphidius colemani</i>	Hymenoptera	0	0	0.25	0	0	0.25	0.32	2
<i>Lysiphlebus testaceipes</i>	Hymenoptera	0	0	0.24	0	0	0.25	0.17	2

# Figure 1

Figure 1: Workflow of the experiment for the three types of pools.

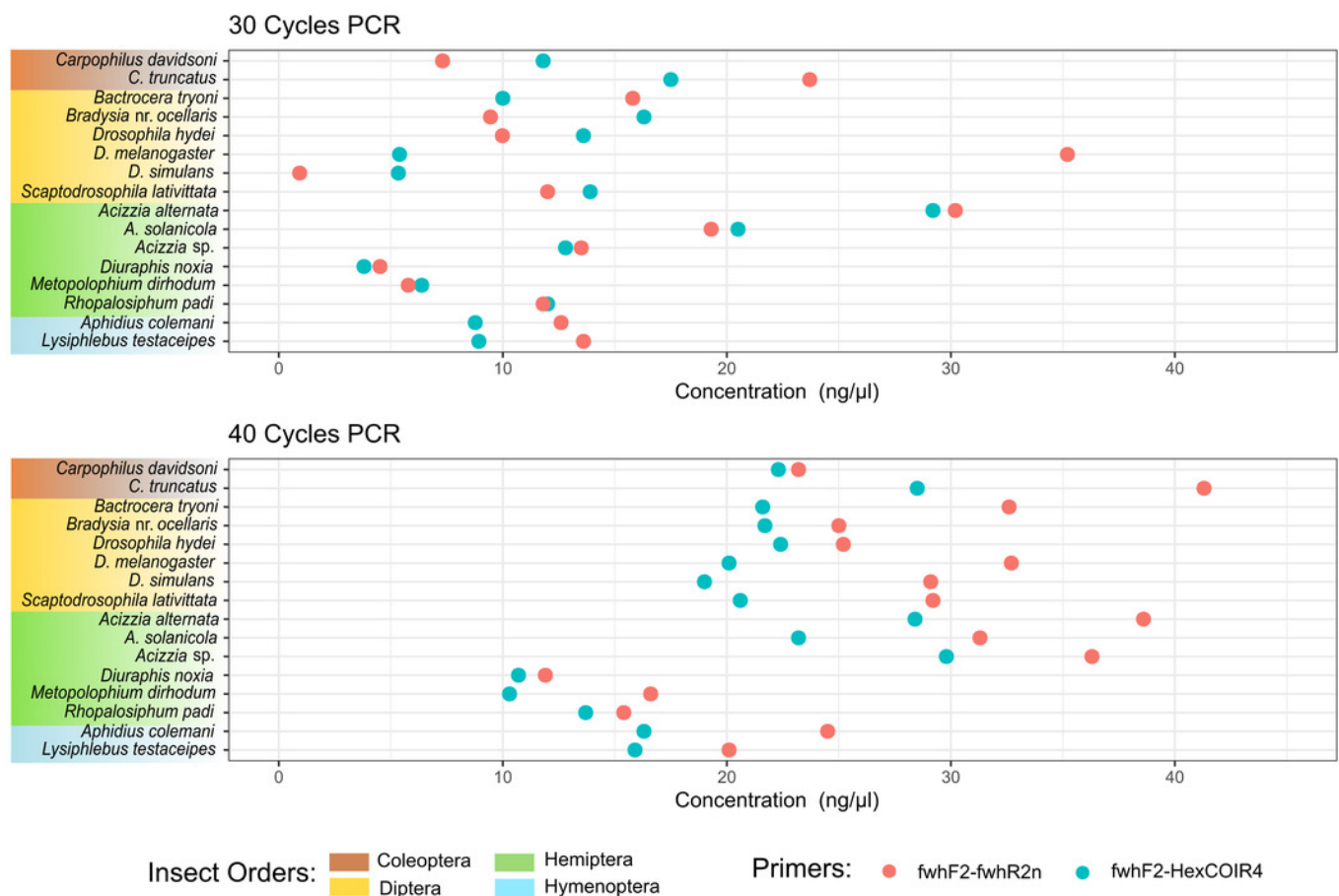
Insect pools are the result of non-destructive DNA extractions from the pooled insect specimens. DNA pools are the result of pooled DNA that was destructively extracted from each insect species. PCR pools are the result of DNA that was extracted, quantified, and amplified separately from each insect species then pooled before indexing qPCR. All pools were indexed, sequenced and analysed following the same protocol. Each pool was amplified using two primer sets.



# Figure 2

Figure 2: DNA concentration of PCR amplification for each species.

DNA concentration was measured from the PCR amplification obtained for each of the two primer pairs after 30 PCR cycles **(A)** or 40 PCR cycles **(B)** on DNA extracted from each single insect species.

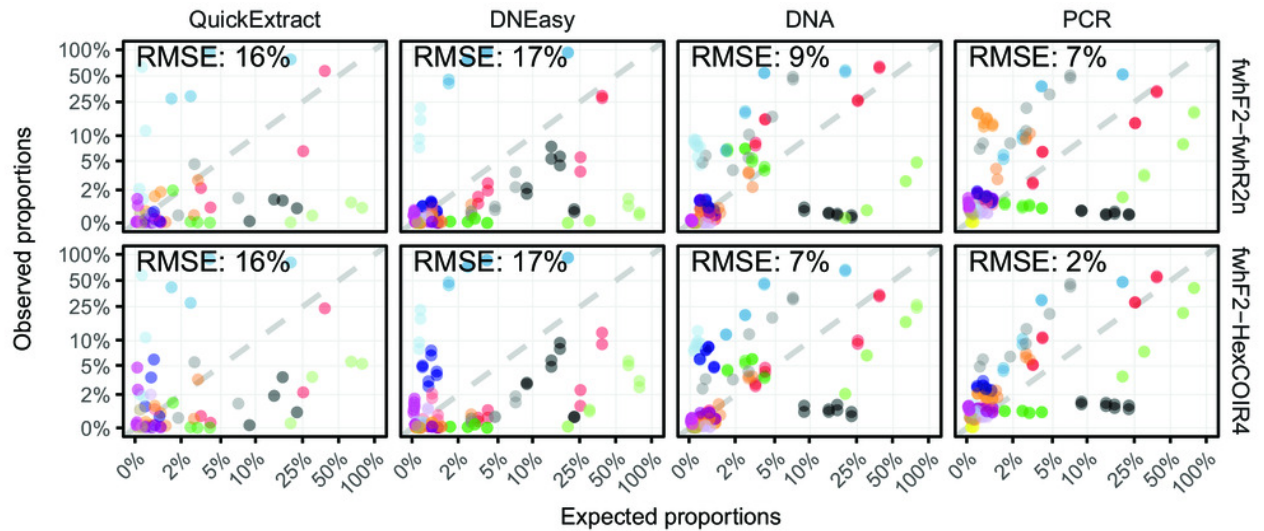


# Figure 3

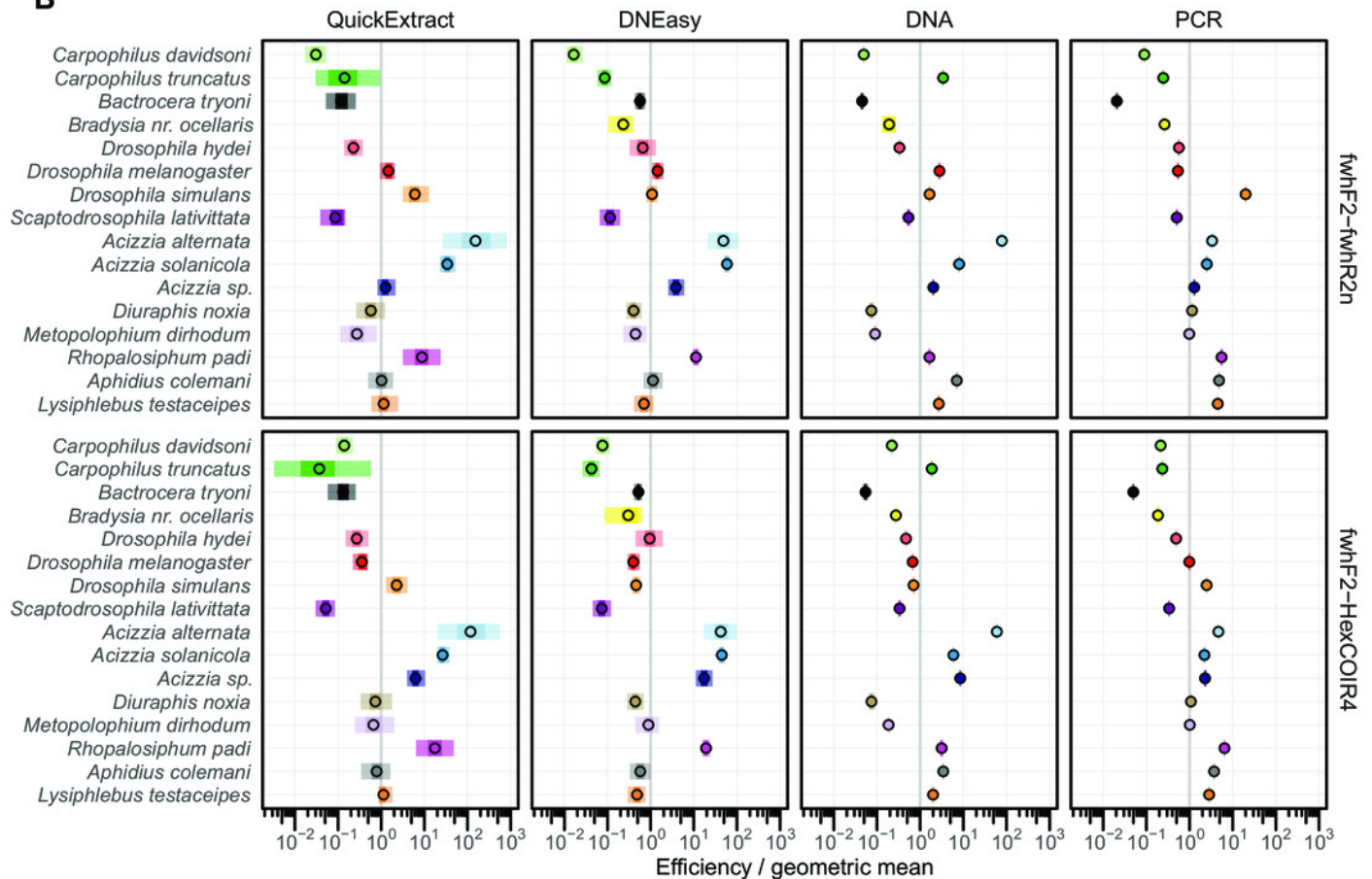
Figure 3: Comparison of estimated and observed proportion of reads.

**A)** Comparison of observed relative abundances to expected relative abundances of each taxon **B)** Estimated bias for each taxon displayed relative to the geometric mean bias, with 95% and 50% confidence intervals displayed. RMSE; Root Mean Square Error.

**A**



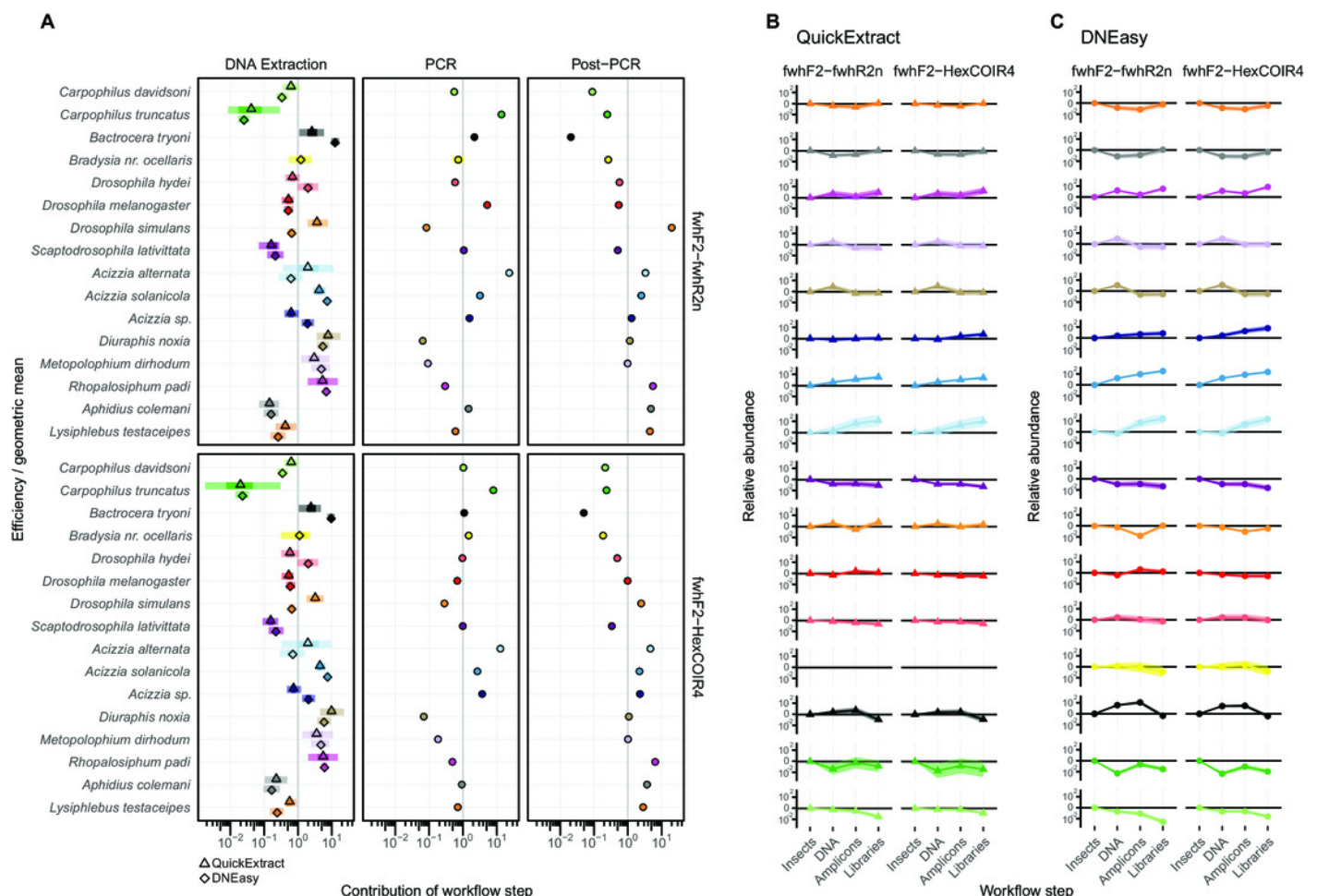
**B**



# Figure 4

Figure 4: Metabarcoding bias partitioned

**A)** Metabarcoding bias partitioned to each protocol and primer step. Change in composition of each taxon throughout the workflow relative to the geometric mean taxon, for **B)** QuickExtract, and **C)** DNEasy obtained by sequentially multiplying a starting abundance of 1 by the partitioned bias estimates from panel A.





# Figure 5

Figure 5: Model predicting the estimated taxon-specific bias.

Coefficients of model predicting the estimated taxon-specific bias from the traits in panel A for the unpartitioned bias and the 3 partitioned bias steps. Coefficients are displayed on a pseudo-log scale to avoid compressing around zero.

