

1 **Title:**

2

3 **MAUI-seq: Metabarcoding using amplicons with unique molecular**
4 **identifiers to improve error correction**

5

6

7 **Authors:**

8 Bryden Fields^{1*} and Sara Moeskjær^{2*}, Ville-Petri Friman¹, Stig U. Andersen², and J. Peter W.
9 Young¹

10

11 *: These authors contributed equally to the work.

12

13 **Author affiliations:**

14 ¹Department of Biology, University of York, York, United Kingdom

15 ²Department of Molecular Biology and Genetics, Aarhus University, Aarhus, Denmark

16

17 **Authors for correspondence:**

18 J. Peter W. Young, peter.young@york.ac.uk

19 Stig U. Andersen, sua@mbg.au.dk

20

21 **Keywords:**

22 Metabarcoding, High-throughput amplicon sequencing, Error correction, Chimeric amplicons,
23 Amplicon sequence variant

1 **Abstract**

2 **Background:** Sequencing and PCR errors are a major challenge when characterising genetic
3 diversity using high-throughput amplicon sequencing (HTAS).

4
5 **Results:** We have developed a multiplexed HTAS method, MAUI-seq, which uses unique
6 molecular identifiers (UMIs) to improve error correction by exploiting variation among
7 sequences associated with a single UMI. We show that two main advantages of this approach
8 are efficient elimination of chimeric and other erroneous reads, outperforming DADA2 and
9 UNOISE3, and the ability to confidently recognise genuine alleles that are present at low
10 abundance or resemble chimeras.

11
12 **Conclusions:** The method provides sensitive and flexible profiling of diversity and is readily
13 adaptable to most HTAS applications, including microbial 16S rRNA profiling and
14 metabarcoding of environmental DNA.

1 Introduction

2
3
4 3 The evaluation of DNA diversity in environmental samples has become a pivotal approach in
5 4 microbial ecology [1] and is increasingly also used to assess the distribution of larger
6 5 organisms [2]. If a core gene can be amplified from environmental DNA with universal primers,
7 6 the relative abundance of species in the community can be estimated from the proportions of
8 7 species-specific variants among the amplicons. High throughput amplicon sequencing
9 8 (HTAS), often termed metabarcoding, has become a cost-effective way to detect multiple
10 9 species simultaneously within a range of environmental samples [3–8]. While shotgun
11 10 sequencing of the whole community (metagenomics) can provide a richer description of the
12 11 functions in a community, HTAS remains a more efficient tool for comparing the species
13 12 diversity of a large number of community samples. Despite the extensive use of HTAS for
14 13 interspecies ecological diversity studies, few investigations have utilised HTAS for
15 14 intraspecies analysis [9, 10]. As 16S rRNA amplicons are too highly conserved to estimate
16 15 microbial within-species diversity, other target gene candidates need to be considered in order
17 16 to sufficiently discern intraspecies sequence variation.

17 17 Many studies have evaluated the extent of PCR-based amplification errors and bias for HTAS
18 18 diversity studies [4, 6, 7, 11]. Numerous known PCR biases reduce the accuracy of diversity
19 19 and abundance estimations, with the major concern being the inability to confidently
20 20 distinguish PCR error from natural sequence variation in environmental samples, which is an
21 21 especially limiting factor for intraspecific studies.

22 22 Polymerase errors, production of chimeric sequences by template switching, and the
23 23 stochasticity of PCR amplification can be major causes of PCR errors [11–13]. Polymerase
24 24 errors introduce new sequences into the template population during amplification. These
25 25 sequence errors include not only substitutions but also insertions and deletions. The use of
26 26 proofreading polymerases, optimised DNA template concentration, and reduced PCR cycle
27 27 number have been suggested to reduce these errors [7, 11, 14].

28 28 In order to account for the introduction of sequence variants in PCR amplification, several
29 29 sequence-classification approaches have been established to manage diversity estimates.
30 30 The most common method is the use of operational taxonomic units (OTUs) in microbial
31 31 diversity studies which analyse target gene sequences and cluster based on an arbitrary fixed
32 32 similarity threshold (QIIME [15]; UPARSE [8, 16–20]). Within species boundaries this technique
33 33 could dramatically reduce the resolution of naturally occurring sequence variation.

34 34 Most recent methods rely on the formation of sequence groups called amplicon sequence
35 35 variants (ASVs) (DADA2, [19]; UNOISE3, [20, 21]). This approach allows sequence resolution
36 36 down to one nucleotide, which is advantageous for determining intraspecies allelic variation,

1 but noise from PCR errors is also more evident. Variation induced by PCR errors often cannot
2 be differentiated from rare natural allelic variation without the use of sequence denoising
3 methods [11]. DADA2 relies on a quality-aware parametric error model, which is developed
4 on a per sequencing run basis. This increases the run time compared to UNOISE3, which
5 uses a one-pass technique [22].

6 An approach that can reduce sequencing noise is to assign a unique molecular identifier (UMI)
7 to every initial DNA template within an HTAS sample, which also enables evaluation of PCR
8 amplification bias [23]. Additionally, the UMI provides a potential route to address polymerase
9 errors in metabarcoding studies. The UMI is provided by a set of random bases in the gene-
10 specific forward inner primer, which introduces a unique DNA sequence into every initial DNA
11 template upstream of the amplicon region during the first round of amplification. Once all
12 original DNA template strands are assigned a unique UMI, an outer forward primer and the
13 gene-specific reverse primer can be used for further amplification. Consequently, all
14 subsequent DNA amplified from the original template will have the same UMI, so the number
15 of reads amplified from the initial template can be calculated. Grouping sequences by shared
16 UMI allows identification of a consensus, which is assumed to be the correct sequence [24].
17 To our knowledge, UMIs have previously only been used for single-amplicon interspecies
18 investigations [25–28].

19 Here, we present a method for metabarcoding using amplicons with unique molecular
20 identifiers to improve error correction – MAUI-seq. The innovative approach is that we use
21 variation among sequences associated with a single UMI to identify erroneous sequences,
22 and we show that this improves error correction compared to non-UMI based analysis using
23 the state-of-the-art software packages DADA2 and UNOISE3.

24 25 **Results**

26 **Laboratory protocol: UMI labelling and amplicon multiplexing**

27 We developed a procedure (MAUI-seq) to amplify multiple target genes from environmental
28 samples, while assigning a random UMI to each initial copy of a template. We opted for a
29 straightforward protocol using a “one-pot” initiation and amplification system. Forward primers
30 consist of two modules; an inner primer bearing the UMI and designed to amplify the target
31 gene, and a universal outer primer that binds only to a linker on the inner primer (**Figure 1A**).
32 We used a 12-base UMI that allowed over 4 million distinct sequences, which is adequate to
33 ensure that duplicate use is negligible for samples with a few thousand sequenced UMIs. For
34 studies with greater sequencing depth, a longer UMI can easily be designed. As a test case,
35 we used MAUI-seq to investigate the genetic diversity of the nitrogen-fixing bacterium
36 *Rhizobium leguminosarum* symbiovar *trifolii* (*Rlt*) by characterising amplicons from the

1 chromosomal core genes *rpoB* and *recA* and the plasmid-borne nodulation genes *nodA* and
2 *nodD*. Each gene was amplified separately in a single reaction, using a target-specific inner
3 forward primer (at low concentration) to assign the UMI and a universal outer primer (at high
4 concentration) to amplify the resulting molecules (**Figure 1A**). The resulting amplicons were
5 pooled and tagged by Nextera to identify the sample, then further pooled for high-throughput
6 paired-end sequencing (**Figure 1B**). The full MAUI-seq step-by-step laboratory protocol can
7 be found in **Additional File 1**.

9 **Analysis protocol: filtering using UMI-based error rates**

10 The resulting paired-end reads were merged and then separated by gene prior to downstream
11 analysis, where UMIs are critical in two ways. Firstly, sequences are clustered by UMI, and
12 the number of unique UMIs is counted for each distinct sequence, selecting the most abundant
13 sequence associated with each UMI (**Figure 1C**). UMIs are discarded as ambiguous if the
14 most abundant sequence does not have at least two reads more than the next in abundance.
15 The most abundant sequence will usually be the correct one (**Figure 2A Case 1**) but, because
16 most UMIs are represented by just a small number of reads, it can sometimes happen that an
17 erroneous sequence is sampled more often than the true sequence, so the primary sequence
18 of the UMI becomes this erroneous sequence (**Figure 2A Case 2**). Secondly, we reasoned
19 that it may be possible to eliminate these errors by using the UMIs to provide information on
20 global error rates across all samples. We implemented this in MAUI-seq by noting both the
21 most abundant (primary) and the second most abundant (secondary) sequence if two or more
22 sequences were associated with the same UMI. MAUI-seq then distinguishes between true
23 and erroneous sequences based on the ratio of primary and secondary occurrences of each
24 sequence, eliminating sequences that show a high ratio (default is 0.7) of secondary to primary
25 occurrences (**Figure 1C** and **Figure 2B**). The 0.7 threshold was chosen empirically, based on
26 the ratios observed for known true and erroneous sequences, but it is a compromise because
27 the incidence of secondary sequences varies across genes and studies. An examination of
28 the results may suggest choosing different thresholds in other studies. Finally, globally rare
29 sequences are discarded (default threshold is 0.1% averaged across samples - a lower
30 threshold could be used if samples were sequenced to a greater depth). Python scripts for
31 separating the genes and for the UMI analysis are available at
32 <https://github.com/jpwyong/MAUI>.

34 **Validation using purified DNA mixed in known proportions**

35 We first evaluated the accuracy of MAUI-seq by profiling DNA mixtures with known strain DNA
36 ratios. DNA was extracted from two *Rlt* strains differing by a minimum of 3bp in each of their
37 *recA*, *rpoB*, *nodA*, and *nodD* amplicon sequences, and the extracted DNA was mixed in

1 different ratios (**Supplementary Table S1**). After amplification and sequencing, assembled
2 reads were assigned to their target gene and analysed using MAUI-seq and two programs
3 frequently used for de-noising of amplicon sequencing data, DADA2 and UNOISE3 [19, 21].
4 Since rare sequences have a high error rate, we discarded (for each of the three methods)
5 sequences that fell below a threshold frequency of 0.1% of accepted sequences. The
6 observed and expected strain ratios were highly correlated for all four genes across the three
7 analysis methods, and we found that the performances of the proofreading (Phusion) and non-
8 proofreading (Platinum) polymerases were gene-dependent, which could be due to
9 differences in amplification efficiency for the four templates (**Table 1** and **Supplementary**
10 **Figures S1-S4**). On average, MAUI-seq detected between 98.5% and 100% true sequences
11 exactly matching those of the two strains in the mixture, while DADA2 ranged from 89.7% to
12 100%, and UNOISE3 from 79.8% to 100% (**Table 1**). The better performance of MAUI-seq
13 was due to more effective elimination of chimeras, which were especially abundant when the
14 PCR reaction was carried out using the Platinum non-proofreading polymerase (**Table 1** and
15 **Supplementary Figures S1-S4**). For the proofreading polymerase, DADA2 detected 100%
16 true sequences for all four genes, whereas MAUI-seq detected 99.03% for *nodA*, failing to
17 eliminate three rare sequences that did not have sufficient secondary counts. This suggests
18 that DADA2 performs equally well or even slightly better than MAUI-seq, when a proofreading
19 polymerase is used to amplify DNA from a simple, two-component mix. The prevalence of
20 secondary sequences varied with gene and polymerase: the secondary/primary ratio for
21 accepted sequences was 0.0322 for *rpoB* using Phusion, but just 0.0002 for *nodD* using
22 Platinum. When the ratio was very low, there were insufficient secondary counts for MAUI-seq
23 to eliminate erroneous sequences effectively.

24 25 **Validation using environmental samples**

26 To test the method on more complex samples, we compared *Rlt* populations in root nodules
27 from two locations in Denmark, a clover trial station in Store Heddinge on Zealand and a lawn
28 at Aarhus University in Jutland (the Field-Samples-1 dataset; **Supplementary Figure S5**).
29 One hundred nodules were pooled for each sample and each plot was sampled in four
30 replicates. Platinum Taq polymerase enzyme was used for amplification. Each clover root
31 nodule is usually colonised by a single *Rhizobium* strain, so a maximum of 100 unique
32 sequences per gene is expected per sample.
33 For Field-Samples-1, the total number of distinct sequences for MAUI-seq and DADA2 were
34 in the same range as the number of distinct alleles observed in a population of 196 natural
35 European *Rlt* isolates [29] (**Table 2**). In contrast, UNOISE3 produced a substantially higher
36 number of distinct sequences, suggesting that its default filtering might be too lenient for our
37 data (**Table 2**). The sequences accepted as true by MAUI-seq were nearly all also included in

1 the DADA2 and UNOISE3 outputs (**Figure 3**). On the other hand, DADA2 and UNOISE3 both
2 accepted a number of sequences that were filtered out by MAUI-seq, and many of these were
3 eliminated by MAUI-seq because a high ratio of secondary to primary occurrences strongly
4 suggested that they represent errors and not real sequences (**Figure 3** and **Additional file**
5 **2**). To provide independent evidence as to whether sequences were likely to be genuine, we
6 checked whether they matched (or differed by a single nucleotide from) known sequences in
7 either a reference database of 196 natural European *Rlt* isolates [29], or the NCBI whole-
8 genome shotgun database (**Figure 3**). The great majority of sequences rejected by MAUI-seq
9 did not have exact matches to these known sequences. A few sequences that exactly
10 matched known alleles were included by DADA2 and UNOISE, but not by MAUI-seq. These
11 sequences were not reported by MAUI-seq because their UMI counts were below the
12 abundance threshold, not because the secondary/primary occurrence filter identified them as
13 erroneous (**Figure 3**). The count threshold could be lowered to include rarer sequences, if
14 the study required it.

15 The allele frequency distributions were different at Aarhus and Store Heddinge (**Figure 3**),
16 and the two sites were clearly separated by the first principal component in a Principal
17 Component analysis (PCA) for MAUI-seq, DADA2 and UNOISE3 sequences. (**Figure 4** and
18 **Supplementary Figure S6-S8**). The amplicon sequencing has sufficient resolution to
19 characterize geospatial variation in allele frequencies. For example, MAUI-seq, DADA2 and
20 UNOISE3 can all clearly identify several highly abundant sequences from one location that
21 are either absent or present in very low frequency in samples from the other location (**Figure**
22 **3**). To quantify the genetic differentiation between the Aarhus and Store Heddinge sites, we
23 calculated fixation indices (F_{ST}). Considering all four target genes combined, the MAUI-seq
24 output resulted in the highest F_{ST} value followed by DADA2 and UNOISE3 (**Table 2**, **Figure 4**
25 and **Supplementary Figure S9-S11**). For all individual genes, MAUI-seq also produced the
26 highest F_{ST} estimates, and the differences were especially pronounced for *nodA*, which also
27 showed the highest overall level of differentiation (**Table 2** and **Supplementary Figure S9-**
28 **S11**). The lower genetic differentiation estimated based on DADA2 and UNOISE3 results,
29 compared to those of MAUI-seq, reflects the inclusion of an increased number of erroneous
30 sequences, which are less differentiated between the two sampled sites than the real
31 sequences (**Figure 3**).

32 Since it was clear from the DNA mixture experiment that the choice of DNA polymerase could
33 significantly affect error rates, we sampled root nodules from 13 additional clover field plots
34 (the Field-Samples-2 dataset) and amplified each sample (a pool of one hundred root nodules)
35 using Platinum and Phusion polymerases in parallel. For samples amplified using Platinum,
36 MAUI-seq detected fewer sequences than DADA2 and UNOISE3 for the two core genes, but
37 the same number of reference sequences were detected (**Table 3**). DADA2 included two

1 chimeric sequences that were filtered out by MAUI-seq due to a high ratio of secondary to
2 primary occurrences (**Additional File 2**). UNOISE3 detected twice as many sequences as
3 DADA2 and MAUI-seq for the accessory genes, but most of the additional sequences had no
4 associated UMIs and were classified as “other” (**Table 3, Additional File 2**). For samples
5 amplified using Phusion, MAUI-seq and DADA2 detected a similar number of sequences
6 (**Table 3**). All nine UNOISE3 *rpoB* sequences that were not accepted by either MAUI-seq or
7 DADA2 (**Additional File 2**) are putative chimeric sequences with two parental sequences of
8 higher abundance. For *nodA*, MAUI-seq includes three sequences that have a single
9 nucleotide difference from a reference sequence, but all have a good ratio of secondary to
10 primary reads, so we hypothesise that these are true sequences. Some reference or exact
11 blast hit sequences were included by DADA2 but not by MAUI-seq because their abundance
12 was estimated by DADA2 to be above the 0.001 threshold, but MAUI-seq estimated that they
13 were rarer.

14 Both MAUI-seq and DADA2 identify and remove sequences that appear to be errors (base
15 substitutions or chimeras), but they use completely different evidence. As a result, they do not
16 always make the same decision, as illustrated for a small set of representative data in **Table**
17 **4** (the *rpoB* sequences amplified by Phusion). While DADA2 examines the sequences and
18 rejects those that are likely to be generated from more abundant sequences in the sample,
19 MAUI-seq does not use the actual sequence but bases decisions on how frequently a
20 sequence occurs as a secondary sequence with the same UMI as another (primary)
21 sequence. Sequences ranked 5 and 6 (**Table 4**) are both potential chimeras of the more
22 abundant sequences 1-4. Both DADA2 and MAUI-seq reject sequence 6 and accept sequence
23 5. Sequence 6 has a secondary/primary ratio of 103/118, which is above the default threshold
24 of 0.7, so MAUI-seq rejects it as a likely error. On the other hand, the ratio for sequence 5 is
25 71/229. This is well below the threshold, but it is higher than other sequences with a similar
26 primary count, e.g. sequence 9 (15/270). A possible explanation is that some of the reads for
27 sequence 5 are generated as chimeras but others are genuine, since is entirely plausible that
28 new alleles are generated by recombination between existing alleles. To some extent, MAUI-
29 seq compensates for this because it allocates sequence 5 a relatively low count and hence
30 lower ranking (8) than it has in the raw reads or the DADA2 analysis. There are two further
31 sequences, 10 and 29, that are rejected by DADA2 as potential chimeras but accepted by
32 MAUI-seq (**Additional file 2** Field-Samples-2-phusion-rpoB); in both cases they have
33 secondary sequence counts well below the threshold, so MAUI-seq accepts them as genuine.
34 DADA2 included an *rpoB* sequence that does not have any associated UMIs (sequence 41),
35 and appears to be a chimera of two more abundant sequences (sequence 3/4/5 and sequence
36 11) (**Table 4**). MAUI-seq counts UMIs, not individual reads, and the default setting is to require
37 that the primary sequence has at least two more reads than the next most frequent sequence

1 (if any) that has the same UMI. This enriches for genuine sequences, which are generally
2 more abundant than errors, but it means, of course, that the number of counts is much lower
3 than the number of reads. In fact, for this particular set of data, the number of UMIs is orders
4 of magnitude smaller than either the raw reads or the DADA2 count, although still sufficient to
5 provide good estimates of the relative abundance of the sequences that make up the bulk of
6 the population. The main reason for the low UMI count is that the number of reads per UMI
7 was suboptimal in these data for the *rpoB* gene: only 18% of the UMIs had more than one
8 read, and MAUI-seq discards single-read UMIs by default. By contrast, in the equivalent data
9 for the *recA* gene in the same study (**Additional file 2** Field-Samples-2-phusion-recA), 37.5%
10 of UMIs had more than one read, making more effective use of the available sequence reads.
11

12 Discussion

13 We propose a new HTAS method (MAUI-seq) designed to assess genetic diversity
14 within or across species. It uses global UMI-based errors rates to detect potential PCR
15 artefacts such as chimeras and single-base substitutions more robustly than the widely-used
16 ASV clustering methods, DADA2 and UNOISE3. The approach is potentially applicable to any
17 study of amplicon diversity, including community diversity estimates based on 16S rRNA and
18 other metabarcoding surveys using environmental DNA.
19

20 Using UMIs to filter out chimeras and other errors

21 In the MAUI-seq approach, UMIs are used to reduce errors in two distinct ways. Since
22 all reads with the same UMI should, in principle, be derived from the same initial template
23 copy, any variation among them reflects errors. In some implementations, a consensus
24 sequence is calculated [24], but we adopt the simpler approach of accepting the most
25 abundant sequence, which will usually give the same result. Requiring more than one identical
26 read before accepting a UMI creates an important quality filter that greatly reduces the number
27 of rare (and usually erroneous) sequences, but as more reads are required, an increasing
28 number of the original reads are discarded and the number of accepted counts declines. To
29 strike a balance between quantity and quality, we chose to count a sequence provided it had
30 at least two more reads than the next most frequent sequence with the same UMI, but this
31 threshold could be adjusted if, for example, a markedly larger number of reads were available.
32

33 While the most abundant sequence associated with a UMI will usually be the correct
34 one, it will sometimes happen that an erroneous sequence will predominate among the small
35 number of reads actually sequenced, leading to these sequences being included among the
36 recorded counts. These errors can be detected, though, by aggregating information across
the whole set of samples. When a UMI is associated with more than one sequence, the

1 secondary sequences are most often erroneous, so sequences that are relatively more
2 abundant as secondary sequences than as the primary sequences associated with UMIs are
3 likely to be erroneous. We recorded the number of times each sequence was found as the
4 second sequence associated with a UMI, and found empirically that a suitable threshold for
5 accepting sequences as genuine was that they occurred less than 0.7 times as often as
6 secondary sequences as they occurred as primary sequences. This threshold can, however,
7 be adjusted to reflect the error distribution observed in a particular study. We found that this
8 approach was very effective in identifying known errors, particularly chimeras, which were
9 generally the most abundant errors. Chimeras were rejected more effectively by MAUI-seq
10 than by the two established ASV clustering methods, DADA2 and UNOISE3. Both of these
11 rely on *de novo* rejection of sequences that could be constructed as recombinants of other
12 sequences that are more abundant in the sample [13]. This method risks rejecting sequences
13 that appear to be recombinant but are genuine alleles, which may not be uncommon,
14 particularly in intraspecific samples. Our approach, by contrast, uses information on the
15 observed error rates in the data (detected using UMIs) to decide whether a sequence is likely
16 to be genuine, regardless of its actual sequence and relationship to other sequences.
17 Sequences that could be generated as chimeras, or that differ by a single nucleotide from a
18 more abundant sequence, may be accepted as genuine if they are more abundant than
19 expected from their rate of occurrence as minor sequences associated with UMIs. In our study,
20 this approach eliminated many known errors and substantially improved our confidence in the
21 remaining data, providing a powerful additional reason for using UMIs in metabarcoding
22 studies of all kinds. While we found that a simple empirical threshold was effective, we noticed
23 that the proportion of secondary sequences varied markedly across studies and genes,
24 suggesting that an adjustable threshold might give further improvement. A useful future
25 development might be to use the abundance of minor sequences associated with UMIs to
26 generate a statistical model of error processes that would provide a firmer theoretical basis for
27 the classification of sequences.

28 29 **Using UMIs to reduce amplification bias**

30 One motivation for the use of UMIs is to obtain more accurate relative abundance data
31 by eliminating possible sequence-specific bias in the PCR amplification, which may be
32 introduced by variation in polymerase and primer affinity for some DNA templates. Indeed, we
33 observed that the Platinum polymerase preferentially amplified the SM170C *rpoB* allele,
34 whereas the Phusion enzyme did not have this bias (**Table 1** and **Supplementary Figure**
35 **S1A-C**). Allele variant bias was also shown for other target genes, although the ranking of the
36 two enzymes was not always the same (**Table 1** and **Supplementary Figures S1-S4**).
37 However, in our study, the use of UMIs did not correct the allele bias. This suggests that the

1 bias was present in the initial round of copying using the target-specific primer, rather than in
2 the subsequent amplification rounds. For our case study, at least, the choice of polymerase
3 was much more important for accurate relative abundance data than the use of UMIs. The
4 main advantage of UMIs was, rather, the ability to remove most sequencing errors, as
5 discussed in the preceding section.
6
7

7 **Advantages of multiplexing several amplicons**

8 Increasing the number of monitored amplicons to four increased our ability to robustly
9 distinguish samples from two locations (**Figure 3-4** and **Supplementary Figure S6-S11**).
10 Multiplexing could be used in other ways, for example to monitor several organisms in the
11 same environment, or to increase read coverage profiling of single genetic markers such as
12 16S [30]. In addition, there is a technical benefit in sequencing multiple different targets
13 together, because a lack of sequence diversity can cause Illumina base-calling issues [31].
14

15 **Optimization of the protocol**

16 As with any metabarcoding project, the first important step is to design the primers
17 carefully to amplify the entire target community with minimum bias, and we used a large
18 database of known gene sequences to achieve this. Another consideration that is shared with
19 other approaches is the choice of polymerase for PCR. For the samples studied here, with
20 abundant template DNA, the proofreading enzyme was clearly superior in performance,
21 although more costly. On the other hand, this enzyme may provide less robust amplification
22 when the template is weak, as we have observed in another project aimed at rhizobial DNA in
23 soil [32]. The use of UMIs introduces other design considerations. We used twelve random
24 nucleotides (with some constraints), giving over four million potential UMI sequences, which
25 was sufficient for the scale of our studies, but it would be simple to increase the UMI length if
26 greater sequencing depth was planned. In any metabarcoding study, the choice of
27 sequencing depth is, to some degree, made blindly because the diversity of templates is not
28 known in advance, but UMI-based approaches need greater depth because it is UMIs that are
29 counted, not reads, and the aim is to have several reads per UMI. There are many factors
30 that affect the average number of reads per UMI, but our study is encouraging in that, without
31 separate optimization, all of our target genes in all of our samples gave usable data. In fact,
32 the number of reads per UMI were suboptimal in most cases. Given a fixed sequencing effort,
33 reads per UMI could, if necessary, be increased by reducing the concentration of the forward
34 UMI-bearing primer and/or of the sample DNA so that fewer distinct UMIs were initiated. With
35 our parameters, at least two reads are needed before a UMI is counted, and a sufficient
36 fraction of the UMIs need at least four reads so that some will have a secondary sequence as
37 well as the primary sequence (with at least two reads more than the secondary).
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1

2 **Future directions for MAUI-seq**

3 HTAS is a valuable and widely-used approach for the study of microbial community
4 diversity, but handling erroneous sequences introduced by the amplification and sequencing
5 procedures has always been challenging. The use of UMIs allows MAUI-seq to greatly reduce
6 the incidence of errors through two mechanisms. Firstly, the requirement that a UMI is
7 associated with at least two identical reads eliminates many rare sequences that are
8 predominantly erroneous. Secondly, sequences that are frequently generated as errors can
9 be identified and removed because they occur unexpectedly often as minor components
10 associated with UMIs that are assigned to more abundant sequences. These mechanisms are
11 independent of any reference database and can recognise and retain genuine alleles that
12 differ by a single nucleotide or match a potential chimera. This makes MAUI-seq particularly
13 suited to studies of intraspecific variation, where the range of sequence divergence may be
14 limited and not fully known in advance. However, the efficient elimination of erroneous
15 sequences is also important in community studies such as those based on widely-used 16S
16 primers, and MAUI-seq should be readily adaptable to this field. The analysis pipeline is very
17 fast because no sequence alignment or database searching is involved; only the accepted
18 final sequences would need to be characterised by comparison to a reference database.

19 Most HTAS studies report the relative proportions of the taxa in a community, but it
20 would sometimes be valuable to estimate the absolute abundance of the microbes in the
21 environmental sample. UMIs can potentially provide such information, if the initial template
22 copying is carefully controlled so that the total number of distinct UMIs reflects the number of
23 templates [26, 33]. While this would necessitate some additional steps at the start of the
24 experimental protocol, it should still be possible to analyse the resulting sequences using the
25 error-removal approaches provided by MAUI-seq. Alternatively, absolute abundance can be
26 estimated by adding a spike of a known quantity of a recognisable target sequence to the
27 sample before processing [11, 34, 35].

28 The addition of a UMI shortens the maximum length of target sequence that can be
29 read, and the counting of UMIs rather than reads requires a higher depth of sequencing, but
30 these limitations are increasingly unimportant as improvements in sequencing technology lead
31 to increasing length, enabling long-read amplicon sequencing [36, 37], and numbers of reads.
32 As implemented in MAUI-seq, UMIs are very effective in reducing the errors inherent in HTAS,
33 and have the potential to improve the quality of any amplicon-based study of diversity.

1 **Materials and methods**

2 **Preparation of DNA mixtures**

3 Two *Rlt* strains (SM3 and SM170C) were chosen based on their *recA*, *rpoB*, *nodA*, and *nodD*
4 sequence divergence, with a minimum of 3 base pair differences in the amplicon region
5 required for each gene. Strains were grown on Tryptone Yeast agar (28°C, 48hrs). Culture
6 was resuspended in 750ul of the DNeasy Powerlyzer PowerSoil DNA isolation kit (QIAGEN,
7 USA) and DNA was extracted following the manufacturer's instructions. DNA sample
8 concentrations were calculated using QuBit (Thermofisher Scientific Inc., USA). DNA samples
9 of the two strains were diluted to the same concentration and mixed in various ratios
10 **(Supplementary Table S1).**

12 **Preparation of environmental samples**

13 For Field-Samples-1 data, white clover (*Trifolium repens*) root nodules were collected from
14 two locations: Store Heddinge, Denmark (6 plots) and Aarhus University Science Park,
15 Aarhus, Denmark (2 plots) **(Supplementary Figure S4)**. The clover varieties sampled were
16 Klondike (Store Heddinge) and wild white clover, (Aarhus). 100 large pink nodules were
17 collected from 4 points on each plot, making a total of 32 samples. Nodules were stored at -
18 20°C until DNA extraction. Nodule samples were thawed at room temperature and crushed
19 using a sterile homogeniser stick. Crushed nodules were mixed with 750µl Bead Solution from
20 the DNeasy PowerLyzer PowerSoil DNA isolation kit (QIAGEN, USA) and DNA was extracted
21 following the manufacturer's instructions. DNA sample concentrations were measured using
22 a Nanodrop 3300 instrument (Thermofisher Scientific Inc., USA).

23 For Field-Samples-2 data, root nodules were additionally sampled from 13 white clover
24 conventionally-managed field trial plots at Store Heddinge, Denmark (Sample 1A-13A,
25 **Additional File 2**). All plots were sown under the same conditions in 2017. Three to ten clover
26 plants were sampled from one point in each plot and the 100 largest nodules collected.
27 Nodules were stored at -20°C, and DNA was extracted for each sample using the Qiagen
28 DNeasy PowerLyzer PowerSoil DNA isolation kit, as above. Samples were processed
29 independently with Platinum (non-proofreading) and Phusion (proofreading) polymerases to
30 evaluate the method dependency on polymerase choice, as described in the following
31 sections.

33 **PCR and purification**

34 Primer sequences were designed for two *Rlt* housekeeping genes, recombinase A (*recA*) and
35 RNA polymerase B (*rpoB*) and for two *Rlt* specific symbiosis genes, *nodA* and *nodD*
36 **(Additional File 1: Table S1).**

1 The three primers are a target-gene forward inner primer, a universal forward outer primer,
2 and a target-gene reverse primer. The concentration of the inner forward primer was 100-fold
3 lower than the universal forward outer primer and the reverse primer (**Figure 1**) in order to
4 reduce the competitiveness of this primer compared to the outer primer. The inner primer is
5 essential for the first round of amplification, but its participation is undesirable in later rounds
6 as it would assign a new unique UMI to an existing amplicon. The PCR reaction mixture and
7 thermocycler programme are provided (**Additional File 1: Tables S2 and S3**).

8 PCRs were undertaken individually for each primer set using Platinum Taq DNA polymerase
9 (ThermoFisher Scientific Inc., USA) (**Additional File 1: Table S2**) and subsequently pooled
10 and purified using AMPure XP Beads following the manufacturer's instructions (Beckman
11 Coulter, USA). Successful PCR amplification was confirmed by running a 0.5X TBE 2%
12 agarose gel at 90V for 2 hours.

13 For the DNA mixture samples, PCRs were run in triplicate. DNA from single strains was also
14 processed as a control to determine the level of cross contamination between samples. Some
15 samples were also amplified using Phusion High-Fidelity polymerase (ThermoFisher Scientific
16 Inc., USA), to evaluate whether use of a proof-reading polymerase improved the quality of the
17 results using the PCR program described in **Additional File 1: Table S2 and Table S4**.

19 **Nextera indexing for multiplexing and MiSeq sequencing**

20 Samples were indexed for multiplexed sequencing libraries with Nextera XT DNA Library
21 Preparation Kit v2 set A (Illumina, USA) using the Phusion High-Fidelity DNA polymerase
22 (ThermoFisher Scientific Inc., USA). PCR reaction mixture and programme are detailed in
23 **Additional File 1: Tables S6 and S7** Indices were added in unique combinations as specified
24 in the manufacturer's instructions (Illumina, USA).

25 The PCR product was purified on a 0.5X TBE 1.5% agarose gel and extracted with the
26 QIAQuick gel extraction kit (QIAGEN, USA) (expected band length: ~454bp). PCR amplicon
27 concentrations were quantified using GelAnalyzer2010a and normalised to 10nM [38]. A
28 pooled sample was quantified and checked for quality by Bioanalyzer (Agilent, USA) before
29 sequencing using Illumina MiSeq (2x300bp paired end reads) by the University of York
30 Technology Facility. A detailed protocol is available in **Additional File 1**.

32 **Read processing and data analysis**

33 The PEAR assembler was used to merge paired ends [39]. Python scripts were used to
34 separate the merged reads by gene (MAUIsortgenes.py) and to calculate allele frequencies
35 both with and without the use of UMIs (MAUIcount.py). The scripts are available in the GitHub
36 repository <https://github.com/jpwyong/MAUI>. Sequences were clustered by UMI, and the
37 number of unique UMIs was counted for each distinct sequence, provided that sequence had

1 at least two more reads with that UMI than any other sequence. In cases where two or more
2 sequences were associated with the same UMI, the second most abundant sequence was
3 noted, and sequences that occurred more than 0.7 times as often as second sequences than
4 as the main sequence associated with a UMI were filtered out of the results as putative PCR-
5 induced chimeras or other errors. Sequences with primers removed (ignoring UMIs) were also
6 clustered using DADA2 (version 1.8) [19] and UNOISE3 (USEARCH version 11.0.667) [21]
7 with default settings. An overall read frequency filter of 0.1% was applied to DADA2 and
8 UNOISE3 outputs to match MAUI-seq accepted sequences filtering. Scripts used for DADA2,
9 UNOISE3, and figure generation are available in **Additional file 3, 4, and 5**, respectively.
10 Output abundance data were then processed for statistical analysis and figure generation
11 using various R packages (**Additional File 3, 4, and 5**; [40, 41]). Principal components were
12 calculated with the R 'prcomp' package using singular value decomposition to explain the
13 *Rhizobium* diversity and abundance within each sub-plot sample. Differences in allele
14 frequencies between samples were quantified using Bray-Curtis beta-diversity estimation
15 using the R package 'vegdist.' PERMANOVA tests were performed using the R package
16 'adonis'. Empirical Bayes estimator of F_{ST} was calculated using the R package 'FinePop'.
17

18 **Acknowledgements**

19 We thank David Sherlock for his experimental expertise in developing this method, the
20 University of York Technology Facility for sequencing, Simon Kelly for DADA2 expertise,
21 Asger Bachmann, Terry Mun, Maria Izabel A. Cavassim, and Marni Tausen for preliminary
22 data analysis and script development, and DLF for access to their clover field trials. This work
23 was funded by grant no. 4105-00007A from Innovation Fund Denmark (S.U.A.). Initial
24 development of the method was funded by the EU FP7-KBBE project LEGATO (J.P.W.Y).
25

26 **Author contributions**

27 Conceptualization: JPWY; Methodology: JPWY, SUA; Software: BF, SM, JPWY; Validation:
28 BF, SM, JPWY, SUA; Formal analysis: BF, SM, JPWY; Investigation: BF, SM; Resources:
29 JPWY, SUA, VPF; Data curation: BF, SM, JPWY; Writing - original draft: BF, SM; Writing -
30 review and editing: BF, SM, JPWY, SUA, VPF; Visualisation: BF, SM; Supervision: JPWY,
31 SUA, VPF; Project administration: JPWY, SUA; Funding acquisition: JPWY, SUA.
32

33 **Availability of supporting data**

34 Raw Illumina reads are available in the SRA repositories with accession numbers [SRP221010
35 (Synthetic mix and Field-Samples-1) and SRP238323 (Field-Samples-2)]. MAUI-seq scripts
36 are available in the GitHub repository <https://github.com/jpwyong/MAUI>. A detailed protocol

1 for sampling, sample preparation, and read processing is available in **Additional file 1**. Scripts
2 used for DADA2, UNOISE3, and figure generation are available in **Additional file 3, 4, and 5**,
3 respectively. Detailed output sequences for all three methods are available in **Additional file**
4 **2**.

6 **Ethics approval and consent to participate**

7 Not applicable

9 **Consent for publication**

10 Not applicable

12 **Competing interests**

13 The authors declare that they have no competing interests.

15 **Funding**

16 This work was funded by grant no. 4105-00007A from Innovation Fund Denmark (S.U.A.).
17 Initial development of the method was funded by the EU FP7-KBBE project LEGATO
18 (J.P.W.Y).

1 **References**

- 2 1. Birtel J, Walser JC, Pichon S, Bürgmann H, Matthews B. Estimating bacterial diversity for
3 ecological studies: Methods, metrics, and assumptions. PLoS ONE. 2015.
- 4 2. Deiner K, Bik HM, Mächler E, Seymour M, Lacoursière-Roussel A, Altermatt F, et al.
5 Environmental DNA metabarcoding: Transforming how we survey animal and plant
6 communities. Molecular Ecology. 2017.
- 7 3. Fonseca VG. Pitfalls in relative abundance estimation using edna metabarcoding.
8 Molecular Ecology Resources. 2018.
- 9 4. Krehenwinkel H, Kennedy SR, Rueda A, Lam A, Gillespie RG. Scaling up DNA barcoding
10 – Primer sets for simple and cost efficient arthropod systematics by multiplex PCR and
11 Illumina amplicon sequencing. Methods in Ecology and Evolution. 2018.
- 12 5. Tessler M, Neumann JS, Afshinnekoo E, Pineda M, Hersch R, Velho LFM, et al. Large-
13 scale differences in microbial biodiversity discovery between 16S amplicon and shotgun
14 sequencing. Scientific Reports. 2017.
- 15 6. Elbrecht V, Leese F. Can DNA-based ecosystem assessments quantify species
16 abundance? Testing primer bias and biomass-sequence relationships with an innovative
17 metabarcoding protocol. PLoS ONE. 2015.
- 18 7. Gohl D, Gohl DM, MacLean A, Hauge A, Becker A, Walek D, et al. An optimized protocol
19 for high-throughput amplicon-based microbiome profiling. Protocol Exchange. 2016.
- 20 8. Poisot T, Péquin B, Gravel D. High-Throughput Sequencing: A Roadmap Toward
21 Community Ecology. Ecology and Evolution. 2013.
- 22 9. Poirier S, Rué O, Peguilhan R, Coeuret G, Zagorec M, Champomier-Vergès MC, et al.
23 Deciphering intra-species bacterial diversity of meat and seafood spoilage microbiota using
24 gyrB amplicon sequencing: A comparative analysis with 16S rDNA V3-V4 amplicon
25 sequencing. PLoS ONE. 2018.
- 26 10. Kinoti WM, Constable FE, Nancarrow N, Plummer KM, Rodoni B. Generic amplicon
27 deep sequencing to determine llarvirus species diversity in Australian Prunus. Frontiers in
28 Microbiology. 2017.
- 29 11. Keschull JM, Zador AM. Sources of PCR-induced distortions in high-throughput
30 sequencing data sets. Nucleic Acids Research. 2015.
- 31 12. Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R. UCHIME improves sensitivity
32 and speed of chimera detection. Bioinformatics. 2011.
- 33 13. Edgar R. UCHIME2: improved chimera prediction for amplicon sequencing. bioRxiv.
34 2016.

- 1 14. Oliver AK, Brown SP, Callaham MA, Jumpponen A. Polymerase matters: Non-
2 proofreading enzymes inflate fungal community richness estimates by up to 15%. *Fungal*
3 *Ecology*. 2015.
- 4 15. Bokulich NA, Kaehler BD, Rideout JR, Dillon M, Bolyen E, Knight R, et al. Optimizing
5 taxonomic classification of marker-gene amplicon sequences with QIIME 2's q2-feature-
6 classifier plugin. *Microbiome*. 2018.
- 7 16. Edgar RC. UPARSE: Highly accurate OTU sequences from microbial amplicon reads.
8 *Nature Methods*. 2013.
- 9 17. Huse SM, Welch DM, Morrison HG, Sogin ML. Ironing out the wrinkles in the rare
10 biosphere through improved OTU clustering. *Environmental Microbiology*. 2010.
- 11 18. Lindahl BD, Nilsson RH, Tedersoo L, Abarenkov K, Carlsen T, Kjøller R, et al. Fungal
12 community analysis by high-throughput sequencing of amplified markers - a user's guide.
13 *New Phytologist*. 2013.
- 14 19. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2:
15 High-resolution sample inference from Illumina amplicon data. *Nature Methods*. 2016.
- 16 20. Fierer N, Brewer T, Choudoir M. Lumping versus splitting – is it time for microbial
17 ecologists to abandon OTUs? 2017.
- 18 21. Edgar RC. UNOISE2: improved error-correction for Illumina 16S and ITS amplicon
19 sequencing. 2016.
- 20 22. Nearing JT, Douglas GM, Comeau AM, Langille MGI. Denoising the Denoisers: An
21 independent evaluation of microbiome sequence error- correction approaches. *PeerJ*. 2018.
- 22 23. Lundberg DS, Yourstone S, Mieczkowski P, Jones CD, Dangl JL. Practical innovations
23 for high-throughput amplicon sequencing. *Nature Methods*. 2013.
- 24 24. Kou R, Lam H, Duan H, Ye L, Jongkam N, Chen W, et al. Benefits and challenges with
25 applying unique molecular identifiers in next generation sequencing to detect low frequency
26 mutations. *PLoS ONE*. 2016.
- 27 25. Faith JJ, Guruge JL, Charbonneau M, Subramanian S, Seedorf H, Goodman AL, et al.
28 The long-term stability of the human gut microbiota. *Science*. 2013.
- 29 26. Hoshino T, Inagaki F. Application of stochastic labeling with random-sequence barcodes
30 for simultaneous quantification and sequencing of environmental 16S rRNA genes. *PLoS*
31 *ONE*. 2017.
- 32 27. Jabara CB, Jones CD, Roach J, Anderson JA, Swanstrom R. Accurate sampling and
33 deep sequencing of the HIV-1 protease gene using a Primer ID. *Proceedings of the National*
34 *Academy of Sciences of the United States of America*. 2011.
- 35 28. Kinde I, Wu J, Papadopoulos N, Kinzler KW, Vogelstein B. Detection and quantification
36 of rare mutations with massively parallel sequencing. *Proceedings of the National Academy*
37 *of Sciences of the United States of America*. 2011.

- 1 29. Cavassim MIA, Moeskjaer S, Moslemi C, Fields B, Bachmann A, Vilhjalmsson B, et al.
2 The genomic architecture of introgression among sibling species of bacteria. *bioRxiv*. 2019.
3 30. Fuks G, Elgart M, Amir A, Zeisel A, Turnbaugh PJ, Soen Y, et al. Combining 16S rRNA
4 gene variable regions enables high-resolution microbial community profiling. *Microbiome*.
5 2018.
6 31. Krueger F, Andrews SR, Osborne CS. Large scale loss of data in low-diversity illumina
7 sequencing libraries can be recovered by deferred cluster calling. *PLoS ONE*. 2011.
8 32. Boivin S, Lahmidi NA, Sherlock D, Bonhomme M, Dijon D, Heulin-Gotty K, et al. Host-
9 specific competitiveness to form nodules in *Rhizobium leguminosarum* symbiovar *viciae*.
10 *New Phytologist*. 2020.
11 33. Kivioja T, Vähärautio A, Karlsson K, Bonke M, Enge M, Linnarsson S, et al. Counting
12 absolute numbers of molecules using unique molecular identifiers. *Nature Methods*. 2012.
13 34. Edgar RC. UNBIAS: An attempt to correct abundance bias in 16S sequencing, with
14 limited success. *bioRxiv*. 2017.
15 35. Palmer JM, Jusino MA, Banik MT, Lindner DL. Non-biological synthetic spike-in controls
16 and the AMPtk software pipeline improve mycobiome data. *PeerJ*. 2018.
17 36. Kumar V, Vollbrecht T, Chernyshev M, Mohan S, Hanst B, Bavafa N, et al. Long-read
18 amplicon denoising. *bioRxiv*. 2018.
19 37. Karst SM, Ziels RM, Kirkegaard RH, Albertsen M. Enabling high-accuracy long-read
20 amplicon sequences using unique molecular identifiers and Nanopore sequencing. *bioRxiv*.
21 2019.
22 38. Lazar I. Gelanalyzer 2010a: Freeware 1d gel electrophoresis image analysis software.
23 2010.
24 39. Zhang J, Kobert K, Flouri T, Stamatakis A. PEAR: A fast and accurate Illumina Paired-
25 End reAd mergeR. *Bioinformatics*. 2014.
26 40. R Core team. R Core Team. R: A Language and Environment for Statistical Computing.
27 R Foundation for Statistical Computing , Vienna, Austria. ISBN 3-900051-07-0, URL
28 <http://www.R-project.org/>. 2015.
29 41. Wickham H. ggplot 2: Elagant graphics for data analysis. 2016.

30

31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Tables

Table 1. Total number of detected sequences in the synthetic mix samples using MAUI-seq, DADA2 and UNOISE3. The percentage of true sequences is averaged over 23 samples for Platinum (non-proofreading) and 14 samples for Phusion (proofreading).

		Platinum			Phusion			
		MAUI-seq	DADA2	UNOISE3	MAUI-seq	DADA2	UNOISE3	exp. seq*
<i>rpoB</i>	n seq	2	3	4	2	2	2	2
	%true*	100	96.96	93.80	100	100	100	-
	Cor.exp/obs*	0.956	0.977	0.981	0.996	0.999	0.9998	-
	chim.freq*	0	0.07	0.13	0	0	0	-
<i>recA</i>	n seq	2	2	2	2	2	2	2
	%true	100	100	100	100	100	100	-
	Cor.exp/obs	0.984	0.991	0.989	0.948	0.952	0.947	-
	chim.freq	0	0	0	0	0	0	-
<i>nodA</i>	n seq	6	5	4	5	2	4	2
	%true	99.04	89.70	89.93	99.03	100	90.43	-
	Cor.exp/obs	0.985	0.998	0.999	0.989	0.999	0.999	-
	chim.freq	0.10	0.25	0.22	0.04	0	0.16	-
<i>nodD</i>	n seq	7	6	21	3	3	14	3 [†]
	%true	98.49	93.93	90.10	100	100	79.83	-
	Cor.exp/obs	0.998	0.998	0.995	0.990	0.998	0.995	-
	chim.freq	0.05	0.05	0.13	0	0	0.11	-
all	%true-overall*	99.76	93.73	91.93	99.74	100	91.71	-

***n seq** is the total number of sequences occurring across all samples. **%true** is calculated by dividing the number of counts for the true sequences by the total number of counts accepted by the method. **%true-overall** is based on summed counts for all four genes. **Cor.exp/obs** is the Pearson correlation for the observed proportion of SM170C reads versus the expected proportion. **Chim.freq** is the proportion of chimeras compared to total reads at 0.5 expected proportion of sequences. **Exp.seq** is the expected number of detected sequences.

[†] SM170C has a second copy of *nodD* [29].

Table 2. Total number of detected sequence clusters in root nodule samples (Field-Samples-1) using MAUI-seq, DADA2, and UNOISE3 clustering and genetic differentiation between populations.

Gene	Method	Detected sequence clusters*					F_{ST}^{\dagger}
		Total	Reference	Exact BLAST	Single nt	Other	
<i>rpoB</i>	MAUI-seq	12	7	3	1	1	0.032
	DADA2	15	7	3	3	2	0.032
	UNOISE3	30	7	2	7	14	0.012
	Reference	13	-	-	-	-	-
<i>recA</i>	MAUI-seq	8	6	2	-	-	0.110
	DADA2	13	8	2	3	-	0.090
	UNOISE3	14	5	2	2	5	0.028
	Reference	17	-	-	-	-	-
<i>nodA</i>	MAUI-seq	9	8	-	1	-	0.369
	DADA2	18	12	1	1	4	0.191
	UNOISE3	43	13	-	5	25	0.061
	Reference	14	-	-	-	-	-
<i>nodD</i>	MAUI-seq	18	11	1	2	4	0.139
	DADA2	22	11	1	3	7	0.124
	UNOISE3	57	11	1	4	41	0.031
	Reference	16	-	-	-	-	-
All genes	MAUI-seq	47	32	6	4	5	0.139
	DADA2	68	38	7	10	13	0.105
	UNOISE3	144	36	5	18	85	0.032

* Output sequences were classified into **reference** (100% identity in at least 1 of 196 *Rhizobium leguminosarum* symbiovar *trifolii* genomes [29]), **exact BLAST** (100% query coverage and 100% identity against the whole-genome shotgun contigs BLAST database), **single nt** (one nt difference from either reference or exact BLAST match), and **other**.

† The population global F_{ST} (fixation index) is an estimate of genetic differentiation among populations based on relative allele abundance.

Table 3. The effect of polymerase choice. Total number of detected sequence clusters in root nodule samples (Field-Samples-2) amplified using Phusion (proofreading) or Platinum (non-proofreading) polymerases. Sequences were clustered using MAUI-seq, DADA2, and UNOISE3.

Gene		Platinum			Phusion		
		MAUI-seq	DADA2	UNOISE3	MAUI-seq	DADA2	UNOISE3
<i>rpoB</i>	Total	16	24	26	15	15	20
	Reference*	9	9	7	8	9	7
	Exact	3	3	2	3	3	2
	Single nt*	3	7	8	3	2	5
	Other*	1	5	9	1	1	6
<i>recA</i>	Total	9	10	12	8	9	10
	Reference	5	5	4	5	5	4
	Exact	0	1	1	0	1	1
	Single nt	3	3	3	3	2	3
	Other	1	1	4	0	1	2
<i>nodA</i>	Total	18	14	35	17	11	34
	Reference	7	10	8	9	9	9
	Exact	0	1	0	0	0	0
	Single nt	6	1	4	6	1	4
	Other	5	2	22	2	1	21
<i>nodD</i>	Total	20	17	46	27	24	71
	Reference	10	12	12	16	16	15
	Exact	0	0	0	0	0	0
	Single nt	6	3	6	5	4	6
	Other	4	2	28	6	3	50

* Output sequences were classified into **reference** (100% identity in at least 1 of 196 *Rhizobium leguminosarum* symbiovar *trifolii* genomes [29]), **exact BLAST** (100% query coverage and 100% identity against the whole-genome shotgun contigs BLAST database), **single nt** (one nt difference from either reference or exact BLAST match), and **other**.

Table 4. A comparison between DADA2 and MAUI-seq for a subset of the Field-Samples-2 data summarised in Table 3: the *rpoB* sequences from samples amplified by Phusion (proofreading) polymerase. Red cells refer to rejected sequences. Green cells refer to sequences, which are accepted by MAUI-seq, while DADA2 rejects them as potential chimeras. Yellow cells refer to sequences filtered out due to low UMI count by MAUI-seq.

Raw reads		MAUI				DADA2		
Rank	count	rank	UMI primary	UMI	accepted	rank	count	accepted
1	99431	1	7459	197	yes	1	54758	yes
2	86751	2	7067	155	yes	2	48402	yes
3	70318	3	3668	95	yes	3	44412	yes
4	47337	4	1898	106	yes	4	28339	yes
5	13190	8	229	71	yes	5	7854	yes
6	11786	9	118	103	no	none	NA	no
7	10490	5	489	19	yes	6	6009	yes
8	9630	6	362	13	yes	7	5414	yes
9	4738	7	270	15	yes	8	2757	yes
10	4290	12	62	15	yes	none	NA	no
11	3223	11	90	3	yes	9	2041	yes
20	1950	10	96	6	yes	10	981	yes
29	1504	13	42	10	yes	none	NA	no
39	1063	14	35	2	yes	12	618	yes
41	946	none	0	0		11	721	yes
43	826	15	34	0	yes	13	434	yes
51	567	16	22	3	yes	14	341	yes
63	415	24	7	0	(yes)	15	208	yes

1

Figure Legends

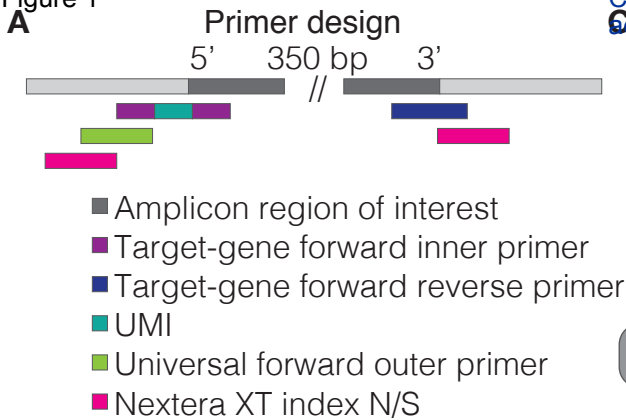
Figure 1. Primer design and method workflow. **A:** Primer design using the sense strand of the target DNA template as an example. The amplicon region of interest should be no longer than 500bp. The target-gene forward inner primer, universal forward outer primer and the target-gene reverse primer are all used in the initial PCR. The Nextera XT indices provide sample barcodes in a separate PCR step. The unique molecular identifier (UMI) region is shown in turquoise on the target-gene forward inner primer. **B:** Sample preparation workflow. **C:** MAUI-seq data analysis workflow.

Figure 2. Erroneous read formation and filtering. **A:** Schematic showing the formation of different sequences with identical UMIs, and bias introduced when sampling for sequencing. **B:** Example data showing the occurrence of real and chimeric *rpoB* sequences as primary and secondary sequence (log scale). S1 and S2: Real sequences derived from two different rhizobium strains (SM170C and SM3). Chi1-4: Chimeric sequences.

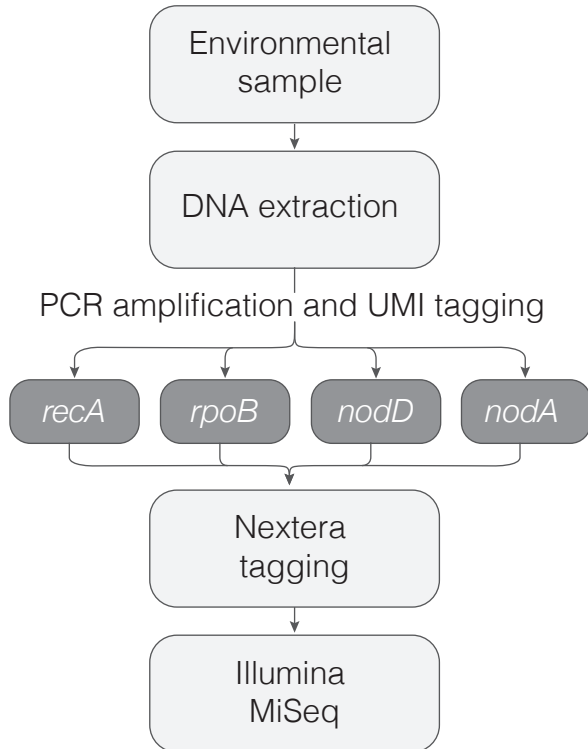
Figure 3. Amplicon diversity reported by MAUI-seq compared with the DADA2 and UNOISE3 analysis pipelines. Data are for four genes from nodule samples from two geographic locations, Store Heddinge (1-6) and Aarhus (7-8). Letters A-D denote the replicates within each plot (**Supplementary Figure 5**). Heatmap of the log₁₀ transformed relative allele abundance of sequence clusters for individual genes. Lines connect identical sequences found by different clustering methods. Evidence that sequences are likely to be genuine is denoted by classifying them as **reference** (100% identity in at least 1 of 196 *Rhizobium leguminosarum* symbiovar *trifolii* genomes [29]), **exact BLAST** (100% query coverage and 100% identity against the whole-genome shotgun contigs BLAST database), **single nt** (one nt difference from either reference or exact BLAST match), and **other**. Sequences not reported by MAUI were classified as **sec/pri ratio** (rejected as erroneous because of a high secondary to primary ratio), **low UMI count** (not reported because too rare), **not found by MAUI** (no accepted UMIs).

Figure 4. Genetic differentiation between populations visualised by Principal Component Analysis (**A-C**) and F_{ST} (**D-F**) of *Rlt* diversity in root nodule samples (8 sites, 4 replicates). Three analysis pipelines are compared: MAUI-seq (**A,D**), DADA2 (**B,E**), UNOISE3 (**C,F**). The PCA analysis was based on log₁₀ transformed relative allele abundance. F_{ST} analysis was based on relative allele abundance. Data from all four genes (*rpoB*, *recA*, *nodA*, and *nodD*) were included in the analysis.

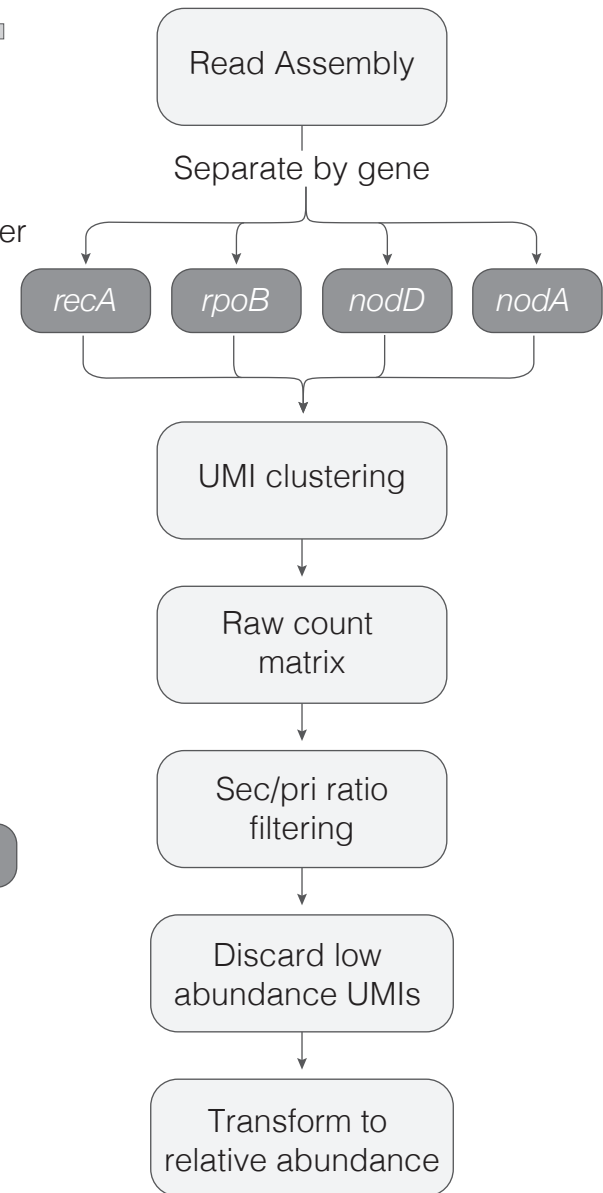
Figure 1

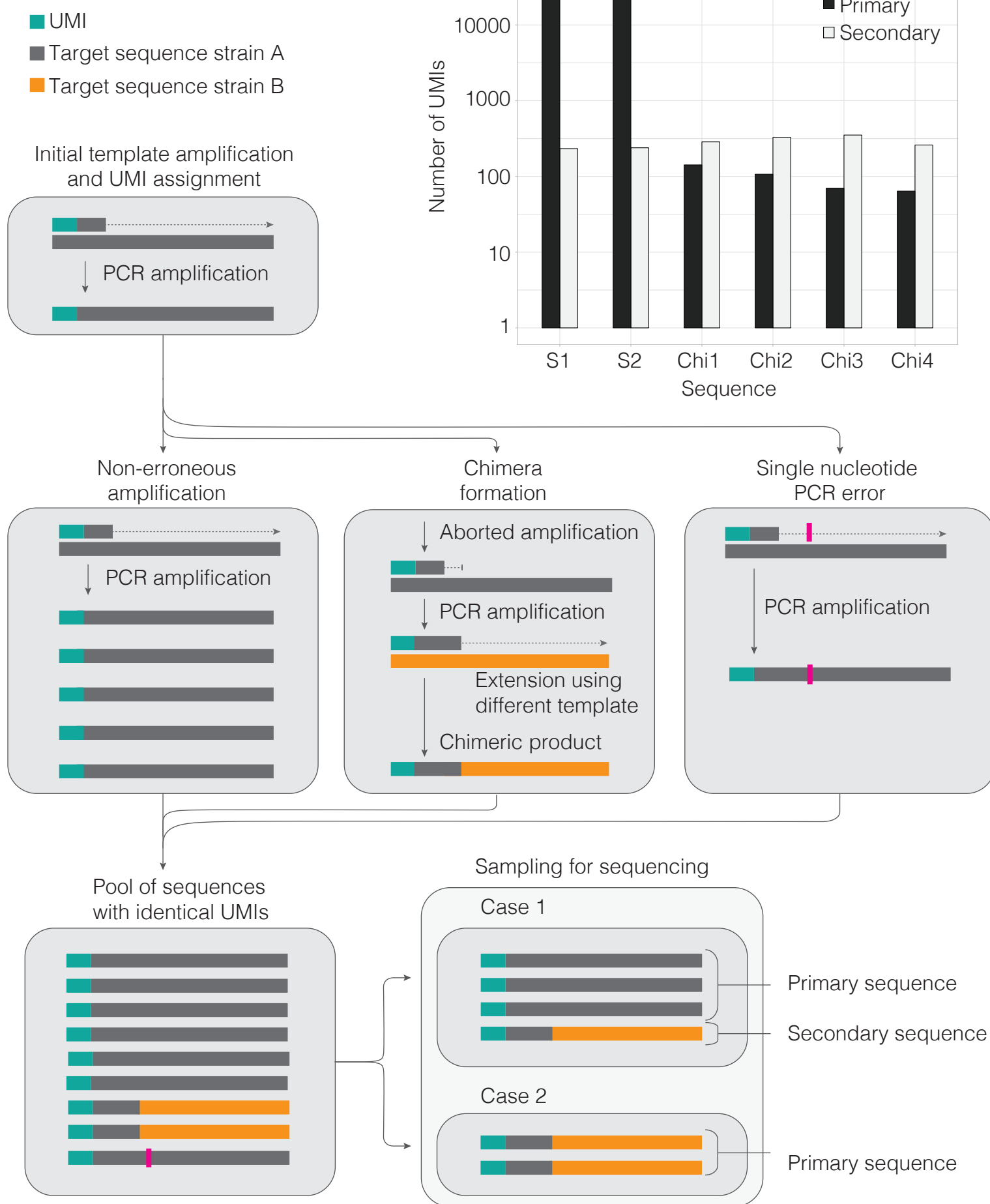


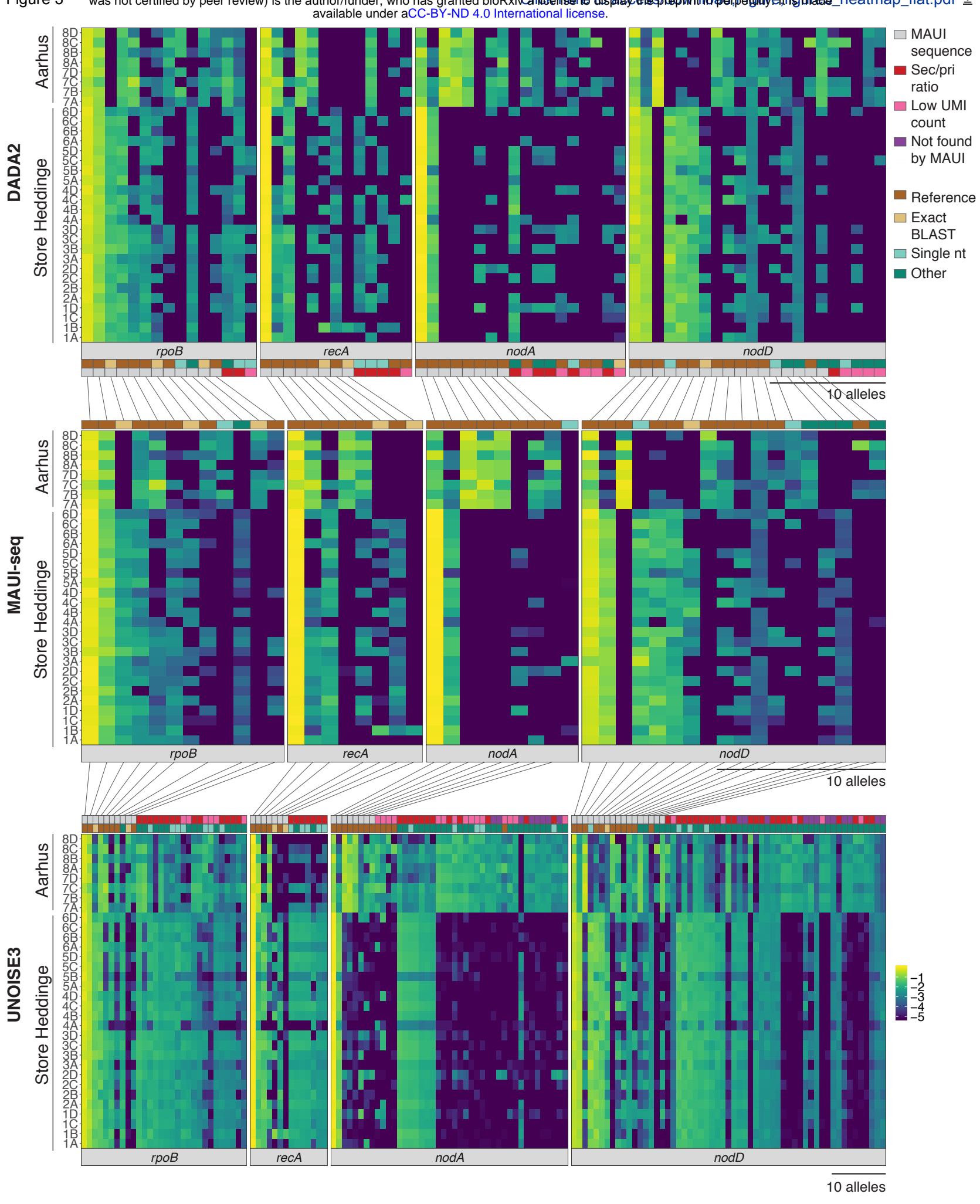
B **Sample preparation workflow**

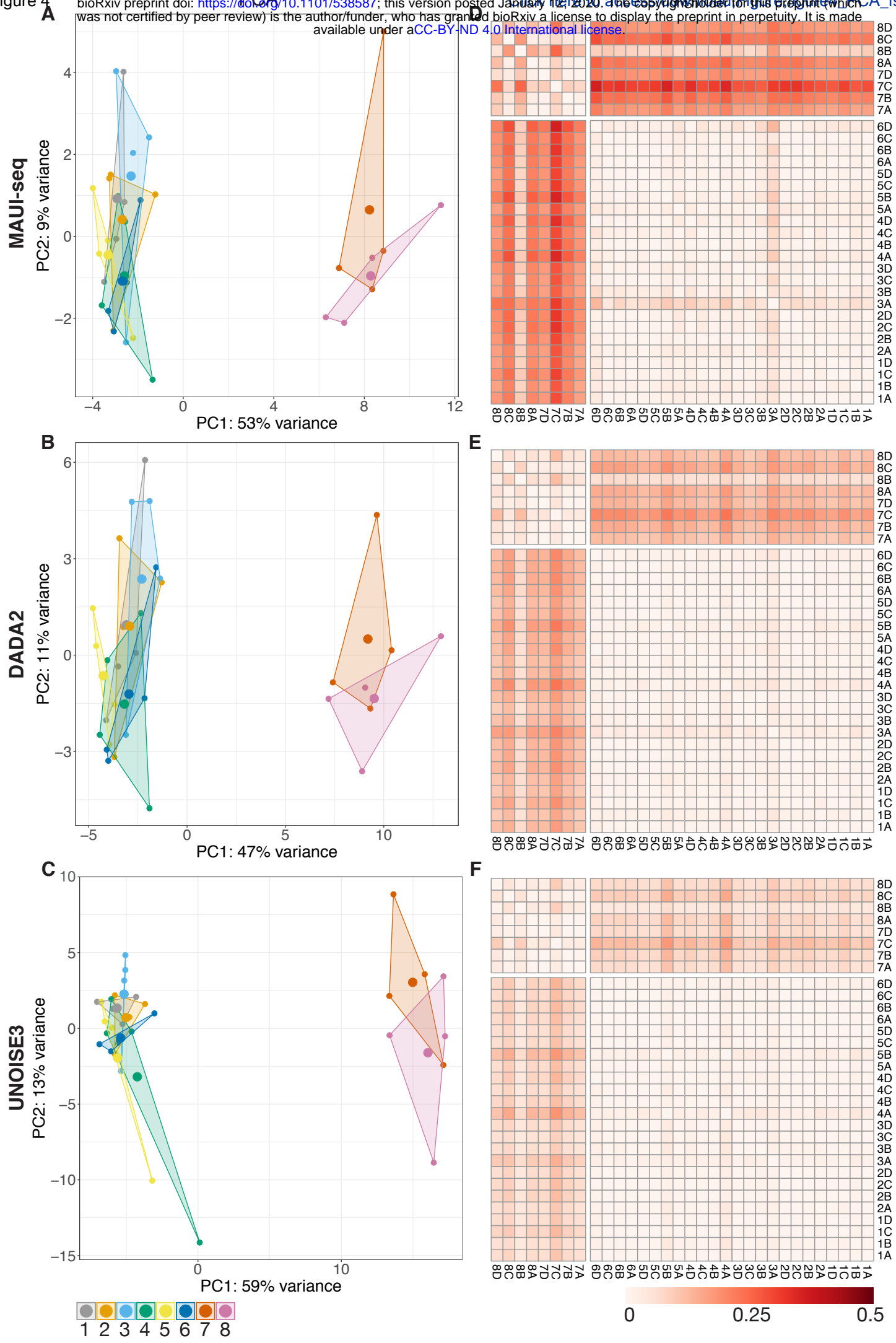


C **Data analysis workflow**

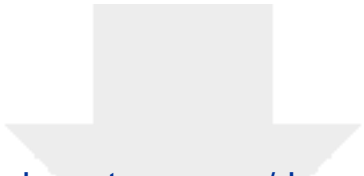







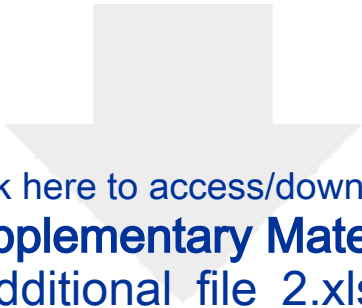







Click here to access/download
Supplementary Material
Additional_file 1.docx

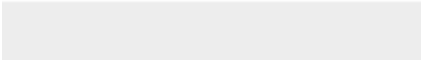





Click here to access/download
Supplementary Material
Additional_file_2.xlsx




Click here to access/download
Supplementary Material
Additional_file_3.R





Click here to access/download
Supplementary Material
Additional_file_4





Click here to access/download
Supplementary Material
Additional_file_5.R

