1  Quantitative evaluation of bias in PCR amplification and Next Generation

2  Sequencing derived from metabarcoding samples

3

4  Marta Pawluczyk1, Julia Weiss1, Matthew G. Links2, Mikel Egaña Aranguren3-4,

5  Mark D. Wilkinson3, Marcos Egea-Cortines1

6  1- Genetics, Instituto de Biotecnología Vegetal, Universidad Politécnica de

7     Cartagena, 30202, Cartagena, Spain

8  2- Department of Computer Science, University of Saskatchewan, Saskatoon

9     Research Centre, 107 Science Place Saskatoon, SK, S7N OX2, Canada

10  3- Centro de Biotecnología y Genómica de Plantas UPM-INIA (CBGP), Campus

11     Montegancedo, Autopista M-40 (Km 38), 28223-Pozuelo de Alarcón Madrid,

12     Spain

13  4- Genomic Resources, Department of Genetics, Physical Anthropology and

14     Animal Physiology, Faculty of Science and Technology, University of Basque

15     Country (UPV/EHU), Sarriena auzoa z/g, 48940 Leioa - Bilbo, Spain

16

17  **Keywords:** meta-barcoding; Next Generation Sequencing; Ion torrent; Ct value;

18  PCR efficiency

19  Corresponding author: Marcos Egea-Cortines, Genetics, Instituto de Biotecnología

20  Vegetal, Universidad Politécnica de Cartagena, 30202, Cartagena, Spain

21  Fax: +34968325433

22  e-mail: marcos.egea@upct.es

1   **Abstract**

2   Unbiased identification of organisms by PCR reactions using universal primers

3   followed by DNA sequencing assumes positive amplification. We used six universal

4   loci spanning 48 plant species and quantified the bias at each step of the identification

5   process from end point PCR to Next-Generation Sequencing. End-point amplification

6   was significantly different for single loci and between species. Quantitative PCR

7   revealed that Ct threshold for various loci, even within a single DNA extraction,

8   showed 2000-fold differences in DNA quantity after amplification. Next Generation

9   Sequencing (NGS) experiments in nine species showed significant biases towards

10  species and specific loci using adaptor-specific primers. NGS sequencing biass may be

11  predicted to some extent by the Ct values of Q-PCR amplification.

12

1

## 1. Introduction

Sequence analysis of complex DNA samples is an important approach to monitoring species distribution in biodiversity and population studies. Genetic material is assessed using universal genomic sequences "barcodes" that are informative regarding the species composition of the sample, as they contain sufficient polymorphisms between species that taxonomic discrimination becomes possible [1]. The barcoding approach has become a mainstream technique to identify species in insects [2], very closely related plant species or hybrids [3], or fungi [4] and bacteria [5].

In plants, seven chloroplast *loci* have been analysed as potential barcodes, the spacers *atpf-atph*, *trnH-psbA*, and *psbK-psbL* , and the genes *matK*, *rbcL*, *rpoB*, *rpoC1* [6, 7]. Metabarcoding involves DNA amplification of barcode loci from mixed population samples, followed by Next-Generation Sequencing (NGS). Sequenced fragments are then either assembled *de novo* and then aligned to known genome sequences [8], or are directly aligned to these genomic databases, thus becoming connected to specific taxa [9]. Most often, the objective of these analyses is to arrive at a quantitative measure of the relative abundance of the various species in the sample.

Despite being a proven tool for taxonomic identification, the approach of PCR is subject to a wide variety of potential biases throughout the processes of amplification and sequence analysis, particularly when applied to mixed-population samples. These biases fall into three main categories. The first relates to differential barcode amplification success as a result of the barcode's universal primers.

3

1    Depending on the marker/species combination, false-negative results can occur when

2    sequence variation at the universal priming sites in one of the species prevents

3    efficient annealing of the universal barcode primer for that species. A second type of

4    bias relates to the efficiency of the amplification reaction, which may differ from

5    species to species based on the sequence composition of their specific variant of the

6    barcode. As a result, the proportion of sequences representing each species in the

7    original sample may bear little resemblance to the proportion of that species in that

8    population. Finally, there may also be biases introduced during the preparation of

9    DNA libraries for sequencing. For instance sample dilution has a strong effect on the

10    correlation between biological and read quantities in bacterial samples [10]. A

11    combination of barcoding and NGS have been in some cases confirmed by qPCR,

12    showing that while the exact quantification is not precise, trends in the population

13    structure are faithful [11].

14    Despite knowing that these potential biases exist, the degree to which each

15    source of bias affects the outcome of a metabarcoding experiment, and their relative

16    importance, have not been well quantified. Moreover, by quantifying these biases and

17    relating them to the specific sequences being studied, it may be possible to formulate

18    approaches for *post facto* normalization of metabarcode data to better-reflect the

19    population make-up. For example, PCR efficiency is an important parameter of

20    Quantitative PCR analysis of gene expression [12–14], and while a variety of

21    algorithms exist that predict the efficiency of PCR amplification, these are currently

22    not considered in any of the normal barcoding or metabarcoding pipelines.

23    Amplification efficiency for a given DNA sequence depends heavily on the G+C

24    content of the amplicon [14]. Under optimal PCR conditions with 100% amplification

25    efficiency, two copies of DNA are generated from each template during exponential

1 phase of amplification, and such a reaction is said to have an efficiency of 2. This

2 efficiency can also affect another important statistic, namely $C_t$, a relative measure of

3 the predicted concentration of the target amplicon in a PCR reaction, and a

4 measurement that is widely used in qPCR analysis [15]. These kinds of statistics will

5 be even more relevant to NGS technologies that introduce additional PCR

6 amplification steps, such as Ion Torrent or 454/Roche that utilize an emulsion PCR

7 during library construction [16].

8 The present study, therefore, aims to first quantitatively analyze PCR success

9 and evaluate amplification efficiency and Ct values as a tool for predicting

10 amplification success. In this study, we undertake a survey of six well-known plant

11 barcoding markers and apply them to 48 species from 34 different plant families. In

12 addition, we apply the Ion Torrent sequencing method simultaneously for mixed

13 species PCR products of three barcoding primers *rbcL*, *rpoB* and *rpoC1* starting with

14 equal amounts of PCR products, to quantitatively measure the bias introduced by this

15 step of the metabarcoding study.

16 Our results reveal that quantitative and even qualitative interpretation of

17 metabarcoding data based on read-abundance is fraught with potential, serious biases.

18 We present, in detail, a dissection of the degree of bias introduced at each step in the

19 typical laboratory practice of barcode marker analysis from mixed DNA samples.

20

21 **2. Materials and Methods**

22 *2.1. Plant material*

1     Plant material 48 plant species belonging to 33 different families was gathered from

2     the local fruit market, field sampling, botanical records and our own collections

3     (Table1).

4     *2.2 DNA extraction and real-time PCR*

5     Two independent genomic DNA samples were extracted from fresh leaf using

6     the commercial kit 'Plant NucleoSpin' (Machery and Nagel, Düren, Germany). All

7     extracted samples were quantified with a Nanodrop and, after isopropanol-ethanol

8     precipitation, diluted to 50 ng/µl. Single species reactions were performed from the

9     two independent DNA extractions with three technical replicas for a total of six PCR

10    reactions per species using 100 ng DNA/reaction. Real-time PCR reactions were

11    performed as described previously [14]. The primers used in this experiment (*rbcL*-a,

12    *matK, rpoB, rpoC1, trnL-F, trnH-psbA*) have been described previously [20] and are

13    presented in Table 2.

14    Equal amounts of genomic DNA from three species were used to create the

15    mixed-species metabarcoding templates. Amplifications were performed using an

16    initial DNA quantity of 150 ng corresponding to 50ng of each of the three genomes.

17    Sequencing reactions comprised nine species.

18

19    *2.3. qPCR efficiency and Ct calculation*

20

21    qPCR efficiency and Ct was computed using *qpcR*, R package [17]. Efficiency

22    value (*E*) was calculated as $E_{cpD2}=F(cpD2)/F(cpD2)-1$, in which F is raw fluorescence

23    at cycle x, and cpD2 is cycle number at second derivative maximum of the curve [18].

1

*2.4. Determination of relative abundance of sequences from PCR products of mixed*

*genomic DNA by semiconductor sequencing*

4

5        PCR products generated by amplifying, separately, the chloroplast barcoding

6 sequences *rbcL*-a, *rpoC1* and *rpoB* from mixed genomic DNAs (100 ng each) were

7 pooled equivalently to yield a final amount of 100ng. Initial time of digestion was

8 adjusted to yield 300 bp fragments. Preparation of samples for library construction

9 and sequencing were performed using the Ion Torrent Next generation sequencing

10 Kits (Life Technologies, CA, USA) according to the manufacturer´s instructions.

11 Briefly PCR products were fragmented using the Ion Shear Plus reagent to a fragment

12 size of 200bp. The corresponding fragments were ligated to adaptors and size

13 fractionated using E-Gel electrophoresis, obtaining fragments of average 330bp.

14 Emulsion PCR was performed using One-touch system according to the manufacturers

15 protocol and sequencing was performed using 314 Ion Torrent chips. A total of

16 333,274 reads with a mean read length of 159bp were computationally analyzed in

17 order to identify species origin of each fragment by aligning the reads with a library of

18 known Chloroplast sequences using Bowtie2 [23]. We extracted from the resulting

19 SAM file a map of reads to the known chloroplast sequences using a Perl script from

20 the mPuma pipeline [8]. The analysis can be reproduced, with the same parameters

21 and data, at the following Galaxy installation. page: http://biordf.org:8983/u/mikel-

22 egana-aranguren/p/sources-of-bias-in-applying-barcoding-markers-for-sequence-

23 analysis-of-environmental-samples.

24

1 **3. Results**

2      This work aimed to reveal and quantify the biases that can occur during

3 metabarcoding analyses. We executed our analyses using the most widely-accepted

4 plant barcodes, quantitated our results using widely-accepted practices such as qPCR,

5 and followed normal protocols for library construction and NGS. At each stage, we

6 re-normalized the samples such that we knew the precise quantities and relative

7 abundances of the input DNA. In addition, although it is known that the size of the

8 PCR amplification product plays a major role in bias within bacterial community

9 pyrosequencing projects [24], the size of the amplicons analysed here is below the

10 1Kb threshold identified in those studies. Thus we should be able to safely exclude

11 that as a possible cause of bias in this study.

12

13 *3.1. Suitability of barcodes depending on plant species*

14      The worst possible outcome of a metabarcode analysis is false-negative, i.e.

15 lack of amplification of a species barcode despite presence of that taxon in the

16 population. As such, our first analysis assessed PCR success. As expected, it varied

17 both between barcode markers, and between the 48 plant species tested. Barcode

18 primers for the *matK* gene were the least successful, giving positive results in only

19 50% of the tested species, followed by *rbcL* which amplified in 82% of species. The

20 *rpoB and rpoC1* genes as well as the short intergenic spacers *trnL –F* and *trnH - psbA*

21 proved to be the most universally successful barcoding markers, amplifying in close to

22 90% of the investigated species. Our data however, gives a within species assessment

23 of PCR success based on six independent amplifications. As none of the samples had a

1    complete failure of amplification with all primer combinations we can conclude that

2    DNA quality was not a limiting factor for amplification.

3

4    *3.2 qPCR parameters for specific barcodes depending on plant species*

5      The second phase of the analysis addressed whether end point PCR results are

6    the outcome of PCR efficiency. As shown in Fig. 2, amplification efficiency during

7    qPCR varied between barcode markers. The highest average efficiency, based on

8    amplification from all species, corresponded to the markers *trnL–F* and *trnH - psbA*

9    followed by *rpoB*, *rpoC1* and *rbcL*. The *matK* barcode showed the lowest average

10    efficiency among all species. The efficiencies of *matK*, *rbcL* and *rpoC1,* but not *rpoB*

11    and *trnH – psbA,* were significantly different from high-efficiency marker *trnL-F*

12    (p<0.0001 for *matK* and *rbc*L and p=0.0013 for *rpoC1*). PCR efficiencies considering

13    all barcode markers for selected species are summarized in Table 3 showing that both

14    the barcode target and the species it is amplified from govern efficiency.

15      Looking at intra-species variation for all barcodes, Ct values varied widely in

16    this case also. Some extreme cases of intraspecific variation were found in *Oryza*

17    *sativa* where *rbcL* showed no amplification whereas *trnL-F* had a Ct of 11.93 (Table

18    3). Beyond the false-negatives, other important differences in Ct were observed for the

19    various markers. In *O. sativa*, the difference in Ct between *matK* (28.55) and *trnL-F*

20    (11.93) is extremely large. If one were to apply the delta-CT formula [15], and

21    assumed an average efficiency for both markers (efficiency = 1.9), the predicted

22    differences in starting DNA level would be 2116-fold based on the estimates from

23    these two barcodes. This was not an isolated case as we found negative amplification

24    of *rbcL* or *matK* and positive albeit differing Ct values in 20% of the species tested for

1    this parameter (*Zea mays, Daucus carota, Quercus coccifera* and *Asphodelus*

2    *fistulosa*).

3        Ct values also varied significantly among species considering all six markers

4    together and these differences did not correlate with the average efficiency of the PCR

5    amplification. For example, *Z. mays* exhibited an average efficiency over all barcodes

6    of 1.88±0.08 and an average Ct of 30.76±4.67, while *Solanum tuberosum* exhibited a

7    similar average efficiency of 1.86±0.15, yet had a Ct of 15.98±5.30.    Moreover, for

8    any given barcode, PCR efficiency and Ct values also proved to be independent

9    variables, based on regression analysis ($R^2$ between 0.37 and 0.003).

10        Differences in efficiency or Ct may be related to amplification bias among

11   template DNAs in environmental samples.    We analysed abundance of reads after

12   sequencing in order to address this question.

13

14   *3.3. Biases during pre-amplification and during emulsion PCR*

15

16        The identification of genomic DNAs corresponding to different organisms in

17   environmental samples requires sequencing of barcode-PCR products. As shown in

18   Fig. 1, not all barcodes successfully amplify in each species. Table 4 shows the result

19   of simultaneous sequencing of equal amounts of PCR products from mixed species

20   templates amplified with barcode markers, *rbcL, rpoB* and *rpoC1*.   The results reveal

21   a strong bias in the number of reads corresponding each species contained in the

22   equimolar  starting  sample.    In  the  case  of  marker  *rpoB*,  most  reads  (95%)

23   corresponded to *Solanum tuberosum* and only 0.02% to *Zea mays*.   The number of

10

1    reads was not related to the PCR efficiencies of the species, but was related to their Ct

2    values when amplified separately (Table 4).

3          Analysis of read numbers also showed a strong bias in the number of total

4    reads corresponding to each of the barcodes (Table 4). Although equal amounts of

5    PCR product from pre-amplification were used to create the amplicon library, only

6    11.2% of all reads were identified as *rbcL* fragments, 36.5% as *rpoB* fragments and

7    52.3% as *rpoC1* fragments. These results are significantly different from an expected

8    33.3% per reaction (Chi-square test p< 2.2 e-16). The relative percentages in read

9    number proved independent of PCR efficiencies of the specific markers but correlated

10    with average Ct values of the marker for the three species amplified.

11          As emulsion PCR for NGS sequencing is performed with primers that

12    correspond to ligated adaptors, and nevertheless a relationship between Ct values and

13    final number of reads is maintained, we can conclude that the main bias that can be

14    encountered in metabarcoding projects is related to the specific sequence of the

15    barcode fragment. This seems to be independent of any primer-specific effect such as

16    internal priming, etc., as it is consistent over two different primer pairs. Library

17    construction can produce at least 4.6 fold differences when comparing *rbcL* against

18    *rpoC1*.

19

20    **4. Discussion**

21          Similarity between primer and template, as well as the regional G+C content of

22    a template, are factors that influence PCR efficiency [19]. The low PCR success,

23    particularly in case of *matK* with 50% PCR failure in a screening of 48 species, is

1    probably due to lack of similarity between primer and template, since no highly

2    conserved sites flanking the most variable parts of this barcoding marker exist [7].

3    The Ct parameter is widely used in Q-PCR analysis [15] and we applied this to

4    assess intraspecific and interspecific variability in both PCR success and as a possible

5    parameter to estimate final read numbers in NGS experiments.  Surprisingly, there was

6    a wide range of Ct values identified within a single species, and even within a single

7    DNA extraction, something completely unexpected as Ct values are thought to relate

8    to DNA/cDNA quantities. These ranges were far beyond the 1-2 cycles that might

9    arise from sampling and manipulation errors.

10    Our results show that PCR efficiency varies among barcoding markers and

11    species, but that these differences in efficiency does not relate to the corresponding Ct

12    values as measure of PCR success. The Ct values in contrast, proved to be a valuable

13    parameter for the estimation of PCR success as *matK* and *rbcL* showed the highest Ct

14    values during qPCR. The late takeoff in the qPCR assay for *rbcL* and *matK* probably

15    reflect an excess of mismatches between primers and templates as Ct values also

16    varied significantly among species over the whole range of markers that may be

17    related to DNA quality and/or PCR inhibiting substances contained in the sample.

18    One of the most common aims in analyzing environmental samples is to

19    estimate the relative abundance of species based on determining the quantity of their

20    template DNAs.  In principle, equal amounts of template DNA from different species

21    should lead to 1:1 amplicon numbers.  However, Suzuki et al. (1996) observed

22    preferential amplification of certain bacterial fragments in mixed templates with lower

23    G+C content [20]. Our results show the situation is similar in plants, with a strong bias

24    in relative read number among three species after Ion Torrent sequencing. Low read

1 numbers corresponded to species with high Ct values for a given marker, whereas

2 PCR efficiency seemed unrelated, indicating that species with lower Ct's for a given

3 marker are preferentially amplified.

4

5 As such, further improving the reliability of amplification, and utilization of sequence

6 content features to derive and apply quantitative data-normalization algorithms, are

7 certainly areas of significant interest for future development in metabarcoding and

8 NGS analysis.

9

10     Acknowledgments

15

16 Literature cited

17 1. Hajibabaei M, Singer GAC, Hebert PDN, Hickey DA (2007) DNA barcoding: how
18     it complements taxonomy, molecular phylogenetics and population genetics.
19     Trends Genet 23:167–172. doi: Doi 10.1016/J.Tig.2007.02.001

20 2. Hajibabaei M, Janzen DH, Burns JM, et al. (2006) DNA barcodes distinguish
21     species of tropical Lepidoptera. Proc Natl Acad Sci U S A 103:968–971.

22 3. Pawluczyk M, Weiss J, Vicente-Colomer MJ, Egea-Cortines M (2012) Two alleles
23     of rpoB and rpoC1 distinguish an endemic European population from Cistus
24     heterophyllus and its putative hybrid ( C. x clausonis) with C. albidus. Plant Syst
25     Evol 298:409–419.

1    4. Krüger M, Stockinger H, Krüger C, et al. (2009) DNA-based species level detection
2       of Glomeromycota: one PCR primer set for all arbuscular mycorrhizal fungi.
3       New Phytol 183:212–23. doi: NPH2835 [pii] 10.1111/j.1469-8137.2009.02835.x

4    5. Links MG, Dumonceaux TJ, Hemmingsen SM, Hill JE (2012) The chaperonin-60
5       universal target is a barcode for bacteria that enables de novo assembly of
6       metagenomic sequence data. PLoS One 7:e49755. doi:
7       10.1371/journal.pone.0049755

8    6. Hollingsworth PM, Forrest LL, Spouge JL, et al. (2009) A DNA barcode for land
9       plants. Proc Natl Acad Sci U S A 106:12794–12797. doi: Doi
10      10.1073/Pnas.0905845106

11    7. Kress WJ, Erickson DL (2007) A two-locus global DNA barcode for land plants:
12      the coding rbcL gene complements the non-coding trnH-psbA spacer region.
13      PLoS One 2:e508. doi: 10.1371/journal.pone.0000508

14    8. Links MG, Chaban B, Hemmingsen SM, et al. (2013) mPUMA: a computational
15      approach to microbiota analysis by de novo assembly of operational taxonomic
16      units based on protein-coding barcode sequences. Microbiome 1:23. doi:
17      10.1186/2049-2618-1-23

18    9. Coissac E, Riaz T, Puillandre N (2012) Bioinformatic challenges for DNA
19      metabarcoding of plants and animals. Mol Ecol 21:1834–47. doi: 10.1111/j.1365-
20      294X.2012.05550.x

21    10. Amend AS, Seifert KA, Bruns TD (2010) Quantifying microbial communities with
22      454 pyrosequencing: does read abundance count? Mol Ecol 19:5555–65. doi:
23      10.1111/j.1365-294X.2010.04898.x

24    11. Links MG, Demeke T, Gräfenhan T, et al. (2014) Simultaneous profiling of seed-
25      associated bacteria and fungi reveals antagonistic interactions between
26      microorganisms within a shared epiphytic microbiome on Triticum and Brassica
27      seeds. New Phytol. doi: 10.1111/nph.12693

28    12. Platts AE, Johnson GD, Linnemann AK, Krawetz SA (2008) Real-time PCR
29      quantification using a variable reaction efficiency model. Anal Biochem
30      380:315–322.

31    13. Pfaffl MW, Horgan GW, Dempfle L (2002) Relative expression software tool
32      (REST(C)) for group-wise comparison and statistical analysis of relative
33      expression results in real-time PCR. Nucl Acids Res 30:e36–. doi:
34      10.1093/nar/30.9.e36

35    14. Mallona I, Weiss J, Egea-Cortines M (2011) pcrEfficiency: a Web tool for PCR
36      amplification efficiency prediction. BMC Bioinformatics 12:404. doi:
37      10.1186/1471-2105-12-404

38    15. Schmittgen TD, Livak KJ (2008) Analyzing real-time PCR data by the
39      comparative CT method. Nat Protoc 3:1101–1108. doi: 10.1038/nprot.2008.73

16. Mardis ER (2008) The impact of next-generation sequencing technology on genetics. Trends Genet 24:133–41. doi: 10.1016/j.tig.2007.12.007

17. Ritz C, Spiess AN (2008) qpcR: an R package for sigmoidal model selection in quantitative real-time polymerase chain reaction analysis. Bioinformatics 24:1549–1551. doi: Doi 10.1093/Bioinformatics/Btn227

18. Spiess A-N, Feig C, Ritz C (2008) Highly accurate sigmoidal fitting of real-time PCR data by introducing a parameter for asymmetry. BMC Bioinformatics 9:221. doi: 10.1186/1471-2105-9-221

19. Polz MF, Cavanaugh CM (1998) Bias in Template-to-Product Ratios in Multitemplate PCR. Appl Envir Microbiol 64:3724–3730.

20. Suzuki M, Giovannoni S (1996) Bias caused by template annealing in the amplification of mixtures of 16S rRNA genes by PCR. Appl Envir Microbiol 62:625–630.

Data availability

Raw and processed data will be made publicly available via entries in Data Dryad, and a formal Data Descriptor will be published detailing the methodologies and workflows used, as well as rich descriptions of the data elements themselves. The analytical workflow for sequence processing and mapping are already publicly available as a Galaxy workflow, as described in the manuscript, and can be freely re-run at any time. The analysis can be reproduced, with the same parameters and data, at the following Galaxy installation. page: http://biordf.org:8983/u/mikel-egana-aranguren/p/sources-

1    of-bias-in-applying-barcoding-markers-for-sequence-analysis-of-environmental-

2    samples.

3

4

5

6

7

1

2    Authors contributions

3    MP, MEC and JW designed experiments, MP and JW performed experiments; MP,

4    JW, MEC, MEA and MDW analyzed data; MP, JW, MEC, MGL and MDW wrote the

5    manuscript. All authors corrected the first draft and approved the manuscript.

6

1    **Table 1** List of plant species analyzed.
2

| Plant species | Family | Location/Donor population |
|---|---|---|
| *Spinacia oleracea* | Amaranthaceae | Murcia, Spain/ commercial |
| *Pistacia lentiscus* | Anacardiaceae | Murcia, Spain/ natural |
| *Daucus carota* | Apiaceae | Murcia, Spain/ commercial |
| *Nerium oleander* | Apocynaceae | Murcia, Spain/ artificial |
| *Arisarum vulgare* | Araceae | Murcia, Spain/ natural |
| *Phoenix dactylifera* | Arecaceae | Murcia, Spain/ commercial |
| *Aloe vera* | Asphodelaceae | Murcia, Spain/ artificial |
| *Lactuca sativa* | Asteraceae | Murcia, Spain/ commercial |
| *Cynara scolymus* | Asteraceae | Murcia, Spain/ commercial |
| *Brassica oleracea botrytis* | Brassicaceae | Murcia, Spain/ commercial |
| *Brassica oleracea italica* | Brassicaceae | Murcia, Spain/ commercial |
| *Diplotaxis erucoides* | Brassicaceae | Murcia, Spain/ natural |
| *Lobularia maritima* | Brassicaceae | Murcia, Spain/ natural |
| *Arabidopsis thaliana* | Brassicaceae | Murcia, Spain/ artificial |
| *Silene vulgaris* | Caryophyllaceae | Murcia, Spain/ natural |
| *Cistus albidus* | Cistaceae | Murcia, Spain/ natural |
| *Cistus heterophyllus* | Cistaceae | Murcia, Spain/ natural |
| *Aeonium arboreum* | Crassulaceae | Murcia, Spain/ natural |
| *Cucumis sativus* | Cucurbitaceae | Biala Podlaska, Poland/ commercial |
| *Ecballium elaterium* | Cucurbitaceae | Murcia, Spain/ natural |
| *Chamaecyparis sp.* | Cupressaceae | Murcia, Spain/ artificial |
| *Arbutus unedo* | Ericaceae | Murcia, Spain/ artificial |
| *Ricinus communis* | Euphorbiaceae | Murcia, Spain/ artificial |

1

| | | |
|---|---|---|
| *Ceratonia siliqua* | Fabaceae | Murcia, Spain/ natural |
| *Pisum sativum* | Fabaceae | Murcia, Spain/ artificial |
| *Vicia faba* | Fabaceae | Murcia, Spain/ artificial |
| *Quercus coccifera* | Fagaceae | Murcia, Spain/ natural |
| *Pelargonium x hortorum* | Geraniaceae | Murcia, Spain/ artificial |
| *Leucobryum glaucum* | Leucobryaceae | Biala Podlaska, Poland/ natural |
| *Anagallis arvensis* | Myrsinaceae | Murcia, Spain/ natural |
| *Callistemos sp.* | Myrtaceae | Murcia, Spain/ artificial |
| *Olea europaea* | Oleaceae | Murcia, Spain/ artificial |
| *Oxalis pes-caprae* | Oxalidaceae | Murcia, Spain/ natural |
| *Pinus silvestres* | Pinaceae | Biala Podlaska, Poland/ natural |
| *Antirrhinum majus* | Plantaginaceae | Murcia, Spain/ artificial |
| *Zea mays* | Poaceae | Murcia, Spain/ commercial |
| *Oryza sativa* | Poaceae | Murcia, Spain/ artificial |
| *Hordeum vulgare* | Poaceae | Murcia, Spain/ commercial |
| *Piptatherum miliaceum* | Poaceae | Murcia, Spain/ natural |
| *Portulacaria afra* | Portulacaceae | Murcia, Spain/ artificial |
| *Galium verrucosum* | Rubiaceae | Murcia, Spain/ natural |
| *Populus alba* | Salicaceae | Murcia, Spain/ artificial |
| *Petunia hybrida* | Solanaceae | Murcia, Spain/ artificial |
| *Solanum tuberosum* | Solenaceae | Murcia, Spain/ commercial |
| *Solanum lycopersicum* | Solenaceae | Murcia, Spain/ commercial |
| *Thymelaea hirsuta* | Thymelaeaceae | Murcia, Spain/ natural |
| *Vitis vinifera* | Vitaceae | Murcia, Spain/ commercial |
| *Asphodelus fistulosus* | Xanthorrhoeaceae | Murcia, Spain/ natural |

1
2
3 1
4 2
5 3
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
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

2  **Table 2** List of primers and amplicon size from the applied barcode markers [7,20]

| DNA region | Primer name | Sequence | Amplicon size (bp) |
|---|---|---|---|
| *rbcL-a* | a_f | ATGTCACCACAAACAGAGACTAAAGC | 670 |
| | a_r | CTTCTGCTACAAATAAGAATCGATCTC | |
| *matK* | 2.1f | CCTATCCATCTGGAAATCTTAG | 857 - 859 |
| | 5r | GTTCTAGCACAAGAAAGTCG | |
| *rpoB* | 2f | ATGCAACGTCAAGCAGTTCC | 548 |
| | 4r | GATCCCAGCATCACAATTCC | |
| *rpoC1* | 1f | GTGGATACACTTCTTGATAATGG | 554 |
| | 3r | TGAGAAAACATAAGTAAACGGGC | |
| *trnH-psbA* | f | ACTGCCTTGATCCACTTGGC | 300 - 389 |
| | f | CGAAGCTCCATCTACAAATGG | |
| *trnL-F* | e | GGTTCAAGTCCCTCTATCCC | 460 |
| | f | ATTTGAACTGGTGACACGAG | |

3

4
5

Table 3. PCR success and qPCR parameters evaluated in a selection of plant species. Samples with NA were non-successful PCR amplifications.

| Plant family | rbcL-a PCR eff | Ct | matK PCRef f | Ct | rpoC1 PCRef f | Ct | rpoB PCRef f | Ct | trnL-F PCRef f | Ct | trnH-psbA PCRef f | Ct | Average ± stdev PCReff | Ct |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Oxalidaceae (*Oxalis pes-caprae*) | 1.89 | 30.99 | 1.83 | 36.24 | 1.70 | 22.63 | 1.78 | 23.44 | 1.91 | 19.41 | 1.90 | 27.76 | **1.84 ± 0.08** | **26.75 ± 6.18** |
| Cistaceae (*Cistus heterophyllus*) | 1.83 | 25.83 | 1.80 | 28.80 | 1.66 | 24.85 | 1.71 | 25.01 | 1.90 | 16.74 | 1.95 | 18.86 | **1.81 ± 0.11** | **23.35 ± 4.58** |
| Poaceae (*Zea mays*) | 1.85 | 34.74 | NA | NA | 1.72 | 22.35 | 1.97 | 25.17 | 1.80 | 20.15 | 1.91 | 26.06 | **1.85 ± 0.10** | **25.69 ± 5.57** |
| Oleaceae (*Olea europaea*) | 1.76 | 26.05 | 1.51 | 23.86 | 1.79 | 17.82 | 1.88 | 15.18 | 1.93 | 16.74 | 1.95 | 17.52 | **1.80 ± 0.16** | **19.53 ± 4.36** |
| Salicaceae (*Populus alba*) | 1.78 | 24.13 | 1.78 | 29.89 | 1.78 | 15.29 | 1.89 | 13.82 | 1.98 | 13.25 | 1.98 | 13.90 | **1.87 ± 0,10** | **18.38 ± 6.96** |
| Poaceae (*Oryza sativa*) | NA | NA | 1.82 | 28.55 | 1.79 | 14.52 | 1.72 | 22.77 | 1.98 | 11.93 | 1.81 | 25.02 | **1.82± 0,10** | **20,56 ± 7.06** |
| Apiaceae (*Dactuca carota*) | 1.94 | 15.82 | NA | NA | 1.85 | 13.06 | 2.00 | 9.77 | 1.98 | 20.15 | 2.00 | 25.95 | **1.95 ± 0.06** | **26.95 ± 6.31** |
| Solananceae (*Solanum tuberosum*) | 1.70 | 16.77 | 1.70 | 20.55 | 1.85 | 10.16 | 1.84 | 8.65 | 1.95 | 10.53 | 2.00 | 10.90 | **1.80 ± 0.12** | **12.93 ± 4.66** |
| Scrophulariaceae (*Antirrhinum majus*) | 1.79 | 27.81 | 1.82 | 33.83 | 1.98 | 13.06 | 1.99 | 12.72 | 2.00 | 12.06 | 2.00 | 15.08 | **1.93 ± 0.1** | **19.09 ± 9.34** |
| Arecaceae (*Phoenix dactylifera*) | 1.87 | 31.39 | 1.90 | 16.06 | 1.97 | 10.81 | 1.97 | 15.32 | 2.00 | 10.12 | 1.84 | 19.95 | **1.92 ± 0.06** | **17.28 ± 7.81** |
| Cucurbitaceae (*Cucumis sativus*) | 1.84 | 27.17 | 1.80 | 29.71 | 1.91 | 9.89 | 1.99 | 9.13 | 1.98 | 9.02 | 1.91 | 23.57 | **1.9 ±0.07** | **18.08 ± 9.77** |
| Amaranthaceae (*Spinacia oleracea*) | 1.90 | 29.66 | 1.42 | 19.59 | 1.99 | 8.94 | 2.00 | 25.32 | 2.00 | 9.40 | 1.99 | 10.40 | **1.88 ± 0.23** | **17.22 ± 8.97** |
| Vitales (*Vitis vinifera*) | 1.82 | 33.15 | 1.85 | 18.17 | 1.75 | 17.65 | 1.94 | 13.66 | 1.89 | 13.88 | 1.95 | 15.48 | **1.87 ± 0.08** | **18.67 ± 7.34** |
| Solanaceae (*Petunia hybrida*) | 1.73 | 28.38 | 1.73 | 19.47 | 1.86 | 11.02 | 1.85 | 10.28 | 1.93 | 10.42 | 1.94 | 11.03 | **1.84 ± 0.09** | **15.10 ± 7.40** |
| Fabaceae (*Ceratonia silique*) | 1.83 | 32.84 | 1.70 | 23.26 | 1.84 | 16.13 | 1.79 | 18.73 | 1.91 | 14.99 | 1.91 | 20.09 | **1.83 ± 0.08** | **21.01 ± 6.50** |
| Fagaceae (*Quercus coccifera*) | NA | NA | NA | NA | 1.68 | 23.39 | 1.72 | 18.43 | 1.90 | 17.06 | 1.86 | 25.14 | **1.79 ± 0.11** | **21.01 ± 3.87** |
| Thymelaeaceae (*Thymelea hirsuta*) | 1.88 | 29.52 | NA | NA | 1.73 | 14.70 | 1.78 | 24.30 | 1.81 | 16.52 | 1.75 | 27.4 | **1.79 ± 0.06** | **22.49 ± 6.58** |
| Xanthorrhoeaceae (*Asphodelus fistulosus*) | 1.81 | 26.73 | NA | NA | 1.73 | 19.38 | 1.76 | 18.13 | 1.78 | 18.91 | 1.84 | 22.84 | **1.78 ± 0.04** | **21.20 ± 3.58** |
| Brasicaceae (*Brassica oleracea*) | 1.70 | 24.55 | NA | NA | 1.76 | 14.76 | 1.82 | 13.57 | 1.76 | 14.35 | 1.67 | 21.83 | **1.74 ± 0.06** | **17.81 ± 5.02** |
| Asteraceae (*Cynara Scolymus*) | 1.49 | 34.47 | 1.62 | 32.27 | 1.50 | 23.89 | 1.49 | 23.45 | 1.49 | 23.27 | 1.40 | 22.94 | **1.5 ± 0.07** | **26.72 ± 5.21** |
| Average | **1.80** | **27.78** | **1.73** | **25.73** | **1.79** | **16.22** | **1.84** | **17.34** | **1.89** | **14.95** | **1.88** | **20.09** | | |
| Stdev | **0.10** | **5.28** | **0.14** | **6.41** | **0.12** | **5.09** | **0.13** | **5.90** | **0.12** | **4.13** | **0.14** | **5.69** | | |

Table 4. Average PCR efficiencies (PCR$_{eff}$,), Ct values and sequence reads derived from PCR products of barcodes *rbcL*, *rpoB* and *rpoC1* using ion semiconductor sequencing

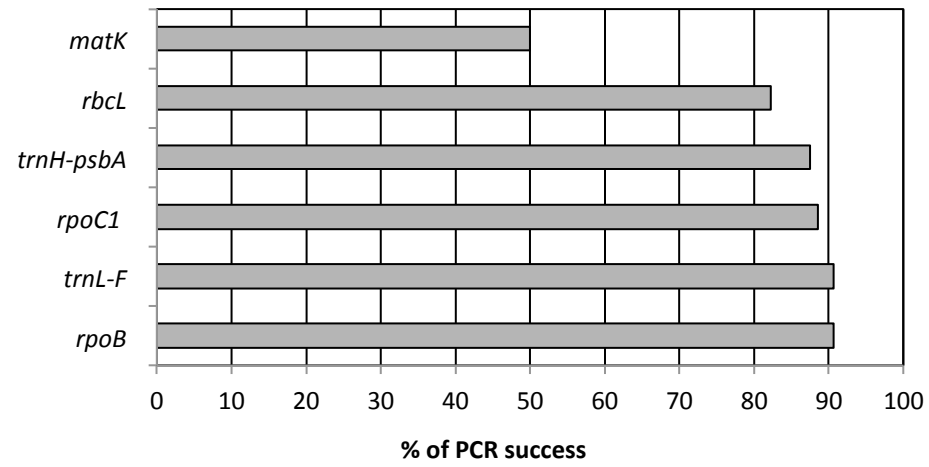| | Barcoding locus | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | *rbcL* | | | *rpoB* | | | *rpoC1* | | |
| Average PCR$_{eff}$ for the amplified species (together) | 1.81±0.09 | | | 1.85±0.14 | | | 1.74±0.06 | | |
| Average Ct for the amplified species (together) | 26.97±7.52 | | | 21.79±5.00 | | | 18.22±4.96 | | |
| Total reads | 34239 | | | 111407 | | | 159923 | | |
| % of total reads | 11.2 | | | 36.5 | | | 52.3 | | |
| | % of species amplified | PCR$_{eff}$ of the species | Ct of the species | % of species amplified | PCR$_{eff}$ of the species | Ct of the species | % of species amplified | PCR$_{eff}$ of the species | Ct of the species |
| *Oxalis pes-caprae* | 0.87 | 1.89±0.04 | 30.99±0.82 | | | | | | |
| *Vitis vinifera* | 4.21 | 1.82±0.02 | 33.15±0.78 | | | | | | |
| *Solanum tuberosum* | 94.92 | 1.69±0.04 | 16.77±0.88 | | | | | | |
| | | | | | | | | | |
| *Zea mays* | | | | 0.02 | 1.71±0.13 | 25.01±0.7 | | | |
| *Cistus heterophyllus* | | | | 1.13 | 1.97±0.06 | 25.17±0.27 | | | |
| *Olea europea* | | | | 98.85 | 1.86±0.01 | 16.28±0.26 | | | |
| | | | | | | | | | |
| *Cistus heterophyllus* | | | | | | | 0.34 | 1.66±0.04 | 24.85±1.24 |
| *Oryza sativa* | | | | | | | 36.57 | 1.79±0.02 | 14.52±0.54 |
| *Populus alba* | | | | | | | 63.09 | 1.78±0.03 | 15.29±1.51 |

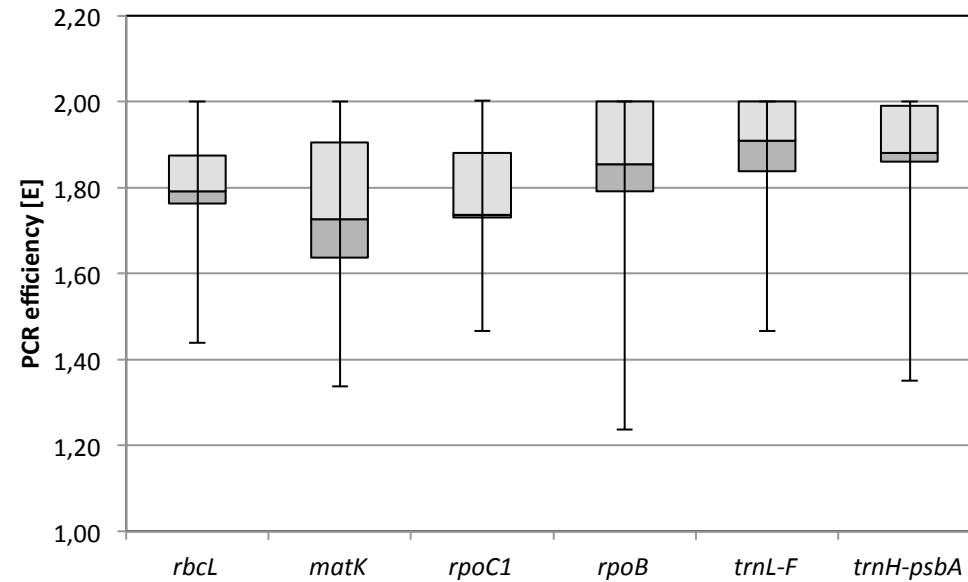Fig. 1. Percent of PCR success of six barcoding markers in a survey of 48 plant species.

Fig. 2. Boxplot of PCR efficiency data for six barcoding markers derived from qPCRs of 48 plant species. The graphic shows only successful amplification data with an efficiency >1.
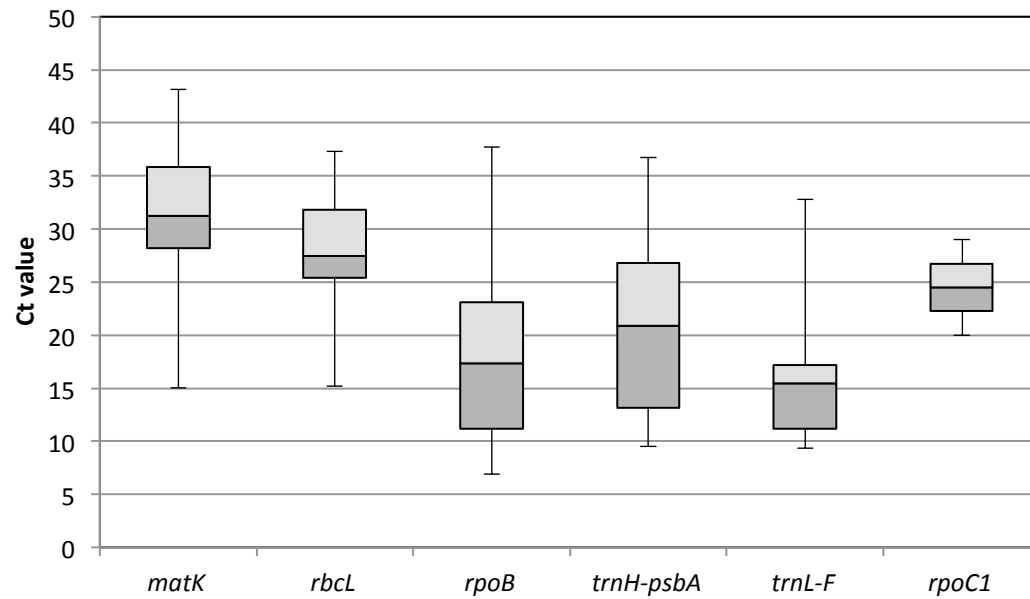
Fig. 3. Boxplot of Ct values for six barcoding markers derived from qPCRs of 48 plant species.