

Bioinformatics analysis on DNA barcode sequences for species identification: a review

ABSTRACT:

Classification of organisms is the primary step in management of biodiversity, breeding, conservation and development of populations and distinguishing adulterant objects. There are many approaches in taxonomic identification, from morphological, PCR-based to sequence-based techniques. Molecular methods give more accurate results than morphological comparisons and are independent of plant stages. PCR-based methods are low-cost but their limited information gives less reproducibility and can only distinguish samples among determined groups. In contrast, in sequence-based methods each nucleotide site is considered as genetic information hence a sequence of nucleotide represents large data, which is highly specific and more stable than PCR bands. Establishment of worldwide DNA library for barcoding is essential. There were previous reviews on screenings and applications of barcodes in different taxa. In this review we discussed common bioinformatics analyses as well as some new improved techniques relying on barcoding approaches.

Key words: molecular classification, bioinformatics tools, DNA barcoding, sequence analysis, identification technique

1. INTRODUCTION

Collection of genetic information, for looking up the origin of a wide range of organisms linked all over the world, is an advanced and essential idea for the protection of species, phylogenetic inference, management and development of genetic diversity [1, 2]. Morphological methods show limitations of accuracy and high reliance on reproductive organs. PCR-based methods can overcome these above problems with just a small piece of sample. However, amplification techniques can only be effectively applied to samples of a defined group using RADP, RFLPs, AFLPs [3] or to samples containing specific genes using species-specific PCR [4-6]. An unknown taxon can not be determined using PCR-band techniques.

Since each site of sequence is considered as a character in bioinformatics analysis, the sequence-based method gives more variable information. This approach allows the read of every nucleotide among the samples. Specific and stable features of monomers are useful in evaluating genetic relationship of a new query sample based on the available sequence library and thus allow to consult origin of the unknown homologous taxon. DNA polymorphism may even provide more information than proteins due to the degradation of genetic code and the presence of large non-coding stretches [7]. DNA fragments represented for the organism can be used as an identifying sequence like the human fingerprint (Figure 1). This is the most common method in molecular identification strategy today. In animal, the highly conserved sequence of mitochondrial cytochrome c oxidase subunit I (COI), which relates to oxidative phosphorylation for metabolism, was effectively used as barcode in diverse animal species [8, 9]. Nevertheless there was no such effective barcode for plants.

A numerous studies of barcoding for plants have been conducted. To date, sequencing techniques for determining species using next generation sequencing (NGS) in plants are limited primarily to agricultural crops [10-13] due to their high cost and time. Therefore short sequences are still considered as convenient and effective tools due to the quick and accurate sequencing [14]. A number of studies have reviewed on finding and applying of DNA barcoding in which the efficiency of different biomarkers have been summarized [15-22]. Others discussed the effectiveness of different barcoding techniques, criteria and measurements [7, 23, 24]. In this review, we discussed some bioinformatics tools in identification analysis of barcoding studies. We also presented some prospects and developed techniques relied on barcoding approaches.

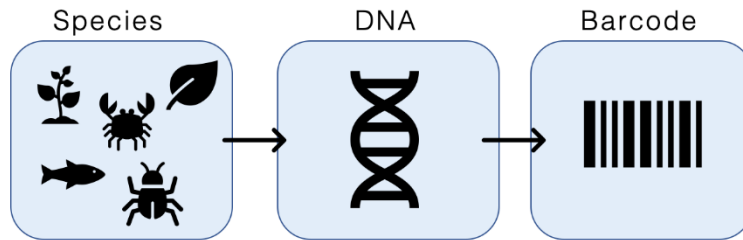


Figure 1. DNA barcoding scheme (https://en.wikipedia.org/wiki/DNA_barcoding)

2. COMMON BIOINFORMATICS ANALYSES IN IDENTIFICATION TECHNIQUES USING SEQUENCES

Bioinformatics procedure for species identification using sequences comprises two basic steps: First, the sequence alignment should be conducted on the basis of a comparative step. This alignment can be performed by two methods. The novel sequences are pairwise aligned against the available sequences in certain databases (similarity search) [25, 26] or studied sequences are aligned against each other in a specific set of data (many-against-each other search) [2, 27, 28]. Following, based on these alignment data, similarity or variation are investigated. This step is performed by evaluating such parameters as GC% content [27, 29], genetic distance [30-33], variable sites [32, 34], indel appearances [35], or monophyletic clusters [36-38], which can indicate typical characters of the examined sequences. These measurements vary from different studies as they depend on alignment and identification methods (Table 1). Final target of this step is to decide whether the species are different or not. In some studies, species resolution was calculated by counting the number of identified species out of the total number of examined species [2, 31, 39].

Table 1. List of common sequence-based identification methods and identification criteria

Alignment Methods	Species Identification Methods	Identification Criteria
Similarity search	BLAST	Correct Identification Ambiguous Identification Incorrect Identification
Many-against-each other search	Genetic distance-based	Nearest distance Correct Identification Ambiguous Identification Incorrect Identification
Many-against-each other search		Best match Correct (match) Ambiguous Incorrect (mismatch)
		Best close match Correct (match) Ambiguous Incorrect (mismatch) Nomatch (under Threshold (%))
		All species barcodes Correct (match) Ambiguous Incorrect (mismatch)
Many-against-each other search		Barcoding gap Intra-specific genetic distance Inter-specific genetic distance Barcoding gap
Many-against-each other search	Tree-based	Monophyletic Paraphyletic Polyphyletic
Many-against-each other search	Character-based (nucleotide polymorphism)	Variable sites Indels Sequence lengths GC% contents

In the alignment step, sequences can be aligned across their entire length (global alignment) or only in certain regions (local alignment). For BLAST (Basic Local Alignment Search Tool)

81 (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>) and FASTA [40] programs, each examined sequence is
82 considered a query sequence. A local pairwise alignment is running between two biological
83 sequences: the query sequence and each of database sequences. In contrast to similarity search,
84 alignment process in many-against-each other search is based on two stages, the pairwise
85 alignments followed by multiple alignments. In multiple alignment, whether global alignment or local
86 alignment should be applied depending on similarity level and difference in sequence lengths.
87 ClustalW is a fully automatic program for global multiple alignment of DNA and protein sequences
88 [41]. However, for MAFFT [42] and MUSCLE [43] programs, users can select suitable aligning
89 strategies. Global alignment is effective when the input sequences share global homology, and the
90 similarity level is high. On the other hand, with fragmentary and divergent sequences, local
91 alignment would be the better option [44].

92
93 For "BLAST" method, the target is searching for the best homologous sequences (best hits) from
94 the available databases (GenBank, BOLD, others). Based on this background, "Correct
95 Identification" means that: the best hit is the sequence with species name as expected. "Ambiguous
96 Identification": the best hits are some sequences belong to different species including the expected
97 species. "Incorrect Identification": the best hit does not match with expected species [25, 26, 28, 45].
98 However, because "BLAST" depends on available sequences reported in databases, query species
99 must be already included in a database otherwise the result will be a failure, therefore a negative
100 "Incorrect Identification" will occur.

101
102 For genetic distance-based methods, the "best hit" feedback of "Correct", "Ambiguous" and
103 "Incorrect" criteria are similar to that in "BLAST", but based on the smallest genetic distances
104 ("Nearest Distance") [28] or the most similarity ("Best match", "Best close match") [46-49]. The
105 query sequence is compared with each other in the given data set. "Best close match" differs from
106 "Best match". A similarity threshold value (e.g. 95%) of all intraspecific distances is established to
107 determine how similar a barcode match needs to be and the results under this similar value (No
108 match) would be removed before identification step.

109
110 For "All species barcodes" methods, a list of sequences sorted by similarity to each query using the
111 same threshold as for "Best close match" are assembled. If the query is followed by all conspecific
112 barcodes with at least two ones then the species is identified. If the query is followed by only one or
113 some of conspecific barcodes then the identification is ambiguous. If the query is followed by none
114 of conspecific barcodes but other species then misidentification occurs [47].

115
116 Barcoding gap method analyses the divergence between intra-specific and inter-specific genetic
117 distances of each query *versus* other conspecific and hetero-specific sequences in a data set [25-
118 28, 30-33, 50, 51]. However this method may be not precise in case of using mean instead of
119 smallest inter-specific distances versus largest intra-specific distances [52]. If barcoding gap exists,
120 the species is successfully identified.

121
122 Another method is the use of nucleotide polymorphism features, such as variable sites, sequence
123 length variation, indel information, GC% content [4, 19, 32, 53, 54] as tools of identification. This
124 approach is known as character-based method in which each nucleotide is consider as the fifth
125 character beside four traditional characters A, T, C and G.

126
127 Regardless of which method is used, the core property of the identification strategy is that all
128 conspecific sequences should be grouped together without blending with any other species. To
129 address this issue, the tree-based method is a simple and visualized approach that is most common
130 in such classification studies [2, 25-27, 29, 34, 36-39, 45, 46, 48, 49, 55]. Two species are
131 completely separated when all sequences of one species are clustered in a monophyletic branch
132 [56]. There are some problems that should be noticed to avoid errors in identification process. Small
133 sample size [4, 26, 29, 36, 51, 57-59] or a wide range but less con-level of taxonomic samples [27,
134 60, 61] may lead to over-fitting [56].

135 136 **3. COMMON BIOINFORMATICS TOOLS IN IDENTIFICATION TECHNIQUES USING** 137 **SEQUENCES**

138
139 For various measurements, different bioinformatics tools could be used also. The BLAST method is
140 presented by BLASTn program from NCBI. Intra- and inter-specific genetic distances, matching
141 sequences, and clustering sequences based on pairwise distances can be calculated in "Best
142 match", "Best close match" and "All species barcodes" methods by Species Identifier tool using

143 TaxonDNA software. When multiple regions are compared for selection of optimal biomarkers,
144 TaxonGap program is used to infer intra- and inter-specific distance for “Nearest” method in high-
145 throughput sequencing researches [62].

146
147 For tree-based reconstruction, probability model should be estimated. Neighbor-joining (NJ),
148 Maximum Likelihood (ML), Maximum Parsimony (MP), Bayesian (BA) or Unweighted Pair Group
149 Method with Arithmetic mean (UPGMA) are common phylogenetic algorithms inferred along with
150 suitable models. Which method should be used in different studies is still a problem that need to be
151 noticed for **obtaining the most accurate results**. Some studies performed comparisons on different
152 methods [27, 45]. Although algorithms are inferred from suppositions, having a thorough
153 understanding of our algorithms and data is the best way to achieve the highest efficiency. The
154 standard principal of tree building is to examine all possible topologies or certain topologies that
155 represent the true structure. Neighbor-joining performance [39, 46] is a time-consuming method in
156 reconstructing phylogenetic trees. It finds pairs of operational taxonomic units (OTUs) that minimize
157 the total branch length at each stage of OTUs clustering and starts with a star-like tree [63].
158 However the reliability of Neighbor-joining tree is a problematic issue [47] as it quickly generates
159 only one phylogenetic tree while others may be more fitting. In contrast, Maximum Likelihood (ML)
160 accounts the probability for all events that can happen simultaneously and the best tree is
161 **supported at a higher probability**. Hence ML become a powerful and professional method in
162 phylogenetic algorithm [64, 65] although it requires significant running time for optimal tree
163 especially with large data [66, 67].

164
165 The Maximum parsimony method is based on the least character state changes required to infer a
166 tree. In case of heterogeneous evolution, maximum parsimony (MP) is strongly biased towards
167 recovering an incorrect tree. However this method outperforms ML over a wide range of
168 conditions, including low and moderate heterogeneity [68].

169
170 While both ML and MP use the probabilities called likelihood, the Bayesian (BA) technique
171 represents the posteriori probability (Bayes’ rule). A known theory called prior is implemented. The
172 posteriori is in direct proportion with the product of likelihood and prior [67]. Bayesian posterior
173 probability gives more-generous estimates of subtree reliability than the maximum likelihood
174 analysis, particularly when using the gamma distribution modeling [65, 67]. When the size of data is
175 small, the probabilities inferred from ML may be over-fitting. Maximum a posteriori (MAP) can be
176 taken accounts to solve this problem.

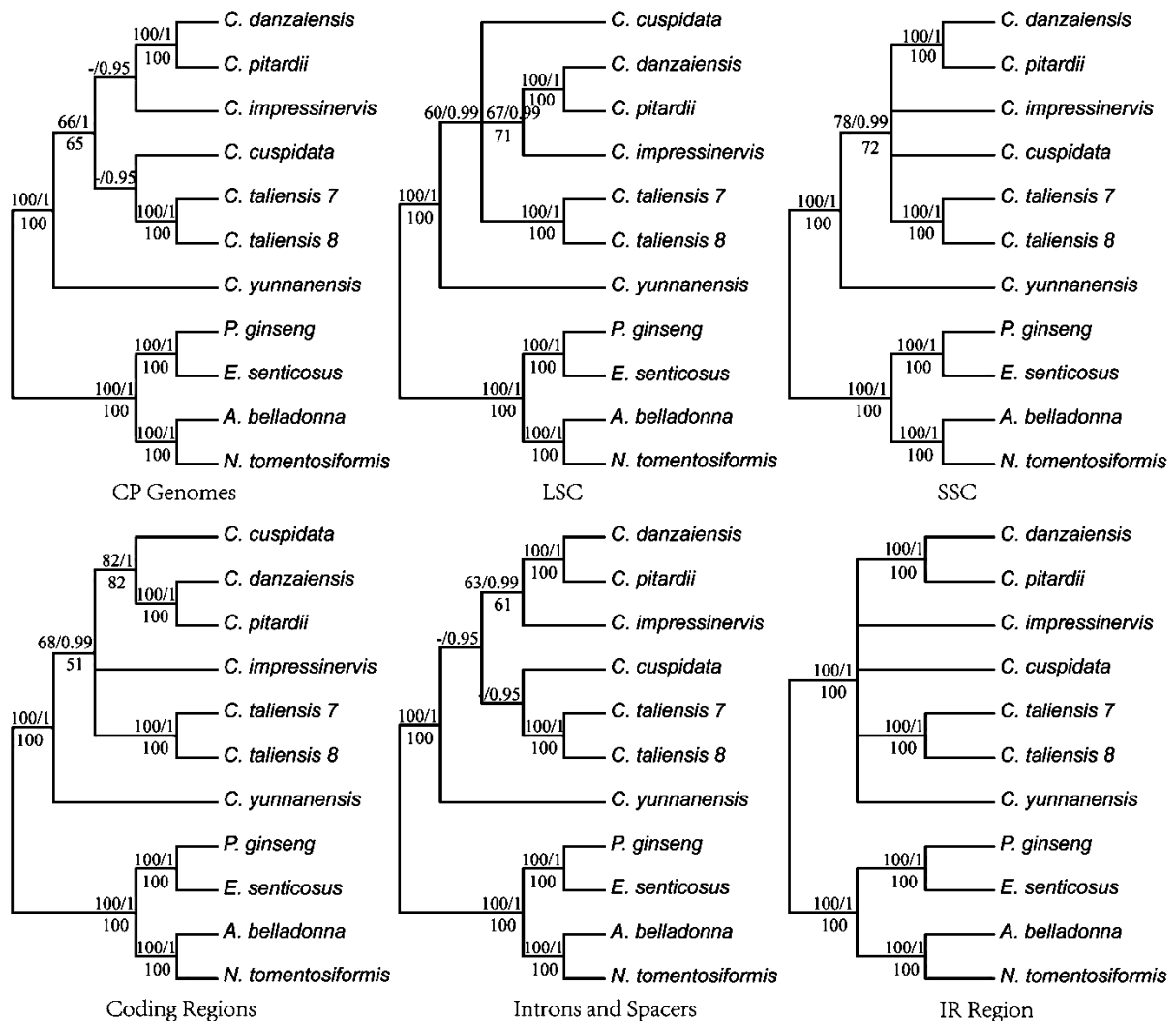
177
178 **Regarding** genetic distance and tree-based methods, MEGA is the most popular software used due
179 to its friendly interface and optimal analysis time. Since genetic variations have to be inferred from
180 evolutionary distance matrices, MEGA versions have integrated these evolutionary models into their
181 program along with different algorithms, i.e. Neighbor-joining, Maximum Likelihood, Maximum
182 Parsimony and UPGMA [69]. However, the number of models in MEGA is pruned for its
183 convenience. For more accurate and reliable results, PAUP* and MRBAYES are often used
184 although it takes much more time. Maximum Likelihood and Maximum Parsimony can also be
185 calculated using PAUP* [70]. Data was estimated for revolutionary models using software
186 jModelTest [71] before running in PAUP*. The software MRBAYES analyses the Bayesian Inference
187 [72]. However, the use of PAUP* and MRBAYES needs bioinformatics skill to run the program,
188 which is quite complex for some researchers.

189 190 **4. RELATION BETWEEN MOLECULAR IDENTIFICATION AND APPLICATION FOR** 191 **PHYLOGENETIC STUDY**

192
193 For discrimination of species, tree-based method seems to be the most common technique in many
194 different studies on variety taxa. Phylogenetic tree can be presented for both phylogenetic and
195 barcoding studies. However phylogenetic analysis represents the measurement and estimation of
196 evolutionary past. Whereas barcoding analysis is used to identify taxa in certain taxonomic groups
197 [2, 39, 48] or to determine a new taxon [60, 73] by DNA sequence comparison. In barcode-
198 phylogenetic tree, **the differences between nucleotide characters are more important than the way**
199 **they form through the evolutionary time**. This means that **taxa of the same species** should be
200 clustered in a monophyletic branch and the different ones should be distributed in separated clades
201 [45, 49]. The length of branches and the members of clusters are not strictly evaluated. Therefore,
202 barcode-phylogenetic tree is not really a phylogenetic tree. However, as the trees are built based on
203 specific DNA sequences, information from them can be used for phylogenetic investigations.
204 Barcoding and phylogenetic relationships of species have also been studied in combination

205 previously [28, 29, 36, 74]. For phylogenetic relationship analyses, homologous variations at each
 206 alignment site are considered. In this case, Maximum Likelihood or Maximum Parsimony **is used**
 207 **instead of** Neighbor-joining [27].

208
 209 Yang *et al.* (2013) used six data sets including sequences of whole complete (cp) genomes, protein-
 210 coding exons, large single-copy region, small single-copy region, inverted repeat region, introns and
 211 spacers for phylogenetic tree establishments to indicate congruent among different data partitions
 212 (Figure 2). Only cp genome tree and the intron and intergenic spacer tree gave genetic similarity.
 213 The relationships between seven species *C. danzaiensis*, *C. pitardii*, *C. impressinervis*, *C.*
 214 *cuspidata*, *C. taliensis* 7, *C. taliensis* 8 and *C. yunnanensis* were not consistent in other trees [75].
 215 This result posed a question: whether a partial DNA sequence region is sufficient to represent
 216 species and the relationship between them. The great data of whole genome might be more reliable
 217 for estimating the evolution compared to shorter barcodes.
 218



219
 220
 221 **Figure 2. Phylogenetic trees constructed from different data partitions from whole**
 222 **chloroplast genome, with all clades were absolutely separated by high genetic variation [75]**
 223 (CP Genome: complete genome; LSC: large single-copy regions; SSC: small single-copy regions; IR
 224 Region: inverted repeat region; numbers above the lines on the left indicate the maximum parsimony
 225 bootstrap of each clade >50%; numbers above the lines on the right indicate the Bayesian posterior
 226 probabilities; numbers below each branch are the maximum likelihood bootstrap of each clade >50%)
 227

228 **5. THE USE OF COMPLETE CHLOROPLAST GENOME AS ULTRA BARCODE**
 229 **(SUPER-BARCODE)**
 230

231 Short DNA sequences **can solve** most questions in identification at species and above-species
 232 levels but still have some limitations with closely related taxa. Some researchers have suggested

233 the way of serving complete chloroplast genome as a single barcode in plants [31, 76]. This method
234 was called ultra-barcoding by Kane *et al.* when they compared nine complete plastid genotypes of
235 *Theobroma cacao* and one related species *T. grandiflorum* with a control GenBank accession [77].
236 The complete chloroplast DNA (cpDNA) could absolutely separate all examined intra-species in the
237 study. This approach promise new effective applications in identification **at species and below-**
238 **species level** [78]. *Fritillaria*, a popular herbal medicine genus in China, experienced some efforts
239 in characterizing closely related species using universal molecular markers (ITS, *trnL-trnF* ect.)
240 but could not be distinguished entirely [79-81]. Using complete chloroplast sequences,
241 phylogenetic analyses of eight *Fritillaria* species were well resolved in the study of Bi *et al.* [82].
242 In other studies, even no potential regions in cp genome were proposed but the complete genome
243 itself had a capability in distinguishing samples as a single barcode [83]. The interesting thing is that
244 you can also use this meta-data to develop potential mini-barcodes which are high variability for
245 quick authentication of certain taxa [84-88].

246
247 This genomic strategy not only allows the use of complete cpDNA as a single barcode, but also
248 facilitate the utilization of other data partitions to identify plants. Protein-coding exons, large single-
249 copy region, small single-copy region, inverted repeat region, introns and intergenic spacers (Figure
250 2) [75], pseudogenes [89] or the differences between size, number of annotated genes [86] could be
251 taken into accounts in phylogenetic analyses. Length-variations which based on the differences of
252 indels (insertions and deletions) at certain locations of cpDNA genome were also considered as
253 effective barcodes [35, 90].

254
255 Genetic information of this super-data is definitely great, enough to avoid the analytical limitation by
256 different bioinformatics methods such as Maximum Likelihood (ML), Maximum Parsimony (MP) or
257 Bayesian (BA). This means that the topologies of phylogenetic trees based on these three
258 algorithms, which be usually congruent using short DNA sequences, are highly similar in the case
259 [75].

260
261 Since the cost for whole genome sequencing has significantly decreased, from \$2.7 billion in 2003
262 for the first human genome to \$300,000 in 2006 and from there to \$1,000 in 2016, a series of
263 studies on plant barcoding was published in the next two years 2017 and 2018. Along with
264 sequencing and assembling techniques, whole-plastome barcode may offer more informative sites
265 and is considered as accurate and effective single barcode for identification in plants.

266 267 **6. APPLICATION OF BARCODES TO DEVELOP OTHER IDENTIFICATION** 268 **TECHNIQUES**

269
270 Nucleotide information can be used to develop some species-specific amplification markers for
271 quick and cheap identification of specific subjects. PCR primers derived from these techniques only
272 react upon annealing to specific DNA sequences and give more specific and reproducible results
273 than random amplifications. SCAR (sequence characterized amplified regions) technique
274 succeeded in authentication three species of *Paphiopedilum armeniacum*, *Paphiopedilum*
275 *micranthum*, *Paphiopedilum delenatii* and their hybrids by developing three species-specific primer
276 pairs from ITS sequences [6]. Some studies focused on developing and comparing simple
277 sequence repeat (SSR) system among screened taxa [86, 91, 92]. Kim *et al.* designed new primers
278 for ARMS (amplification refractory mutation system) technique based on specific insertion in
279 sequence of *Cypripedium macranthos* and SNPs (single nucleotide polymorphisms) in sequences
280 of *Cypripedium japonicum* and *Cypripedium formosanum*) located inside *atpF-atpH* barcode. These
281 three primer pairs can be used in combination to distinguish four *Cypripedium* species with different-
282 size bands on the electrophoresis gel [4]. RFLP (restriction fragment length polymorphism)
283 technique is a type of random PCR-based method in identification of subjects. However, RFLP
284 based on the combination of ITS and some chloroplast sequences gave more reproducible and
285 successful identification of native *Dendrobium* species in Thailand by Peyachoknagula *et al.* [93].
286 Therefore using known barcodes to develop other **time- and cost-saving methods** can also support
287 molecular identification.

288
289 Furthermore, metabarcoding using high-through put sequencing can help identify a variety of
290 species in multiple samples. A detection of species was performed simultaneously in 55 commercial
291 salep products based on ITS barcode. Each sample was found to contain 1-55 species [94]. *RbcL*,
292 *matK* and ITS were also used in combination to determine 16 orchid species as components of a
293 **common food named** *Chikanda* in Zambia [95]. Based on this foundation, the authors alerted the

294 over-harvesting condition and called for the conservation of these rare orchids. This technique also
295 opened a new trend in application of DNA barcodes.

296

297 **7. CONCLUSIONS**

298

299 Barcoding technique had been shown to be useful in many practical applications and classification
300 studies. This method promises more optimal results than the PCR method especially when the price
301 of the sequencing is decreasing. In barcoding technique using sequences, the chosen loci and
302 algorithms are directly related to the identification results.

303

304 Whole genome comparison based on next generation sequencing and high-throughput sequencing
305 has been shown to be more convenient and contain greater data than traditional barcoding.
306 However, in terms of price, it is about \$200 per sample up to now, which is 20 times more
307 expensive than mini-barcodes (\$11-13). In terms of technology, a powerful computer, which is not
308 available in small laboratories, is needed to manage large genome data. Besides, it also takes a
309 considerable time from the sequencing to the analyzing. For those reasons, traditional barcoding
310 methods are still popular and effective in many cases. Hence these two barcoding trends can be
311 performed according to the demands and conditions of different laboratories. In parallel, the
312 sequencing techniques and analyzing tools should be improved to make it simpler and more
313 convenient for researchers. The next orientation of identification should be a species identification
314 gadget that is portable and requires no amplification step.

315

316 **ACKNOWLEDGEMENT**

317 The study was supported by Air Force Office of Scientific Research and Asian Office of Aerospace
318 Research and Development (grant number FA23861514119). All authors declared no other
319 competing financial interests.

320

321 **ABBREVIATION**

322	AFLP	Amplified Fragment Length Polymorphism
323	ARMS	Amplification Refractory Mutation System
324	BA	Bayesian
325	BLAST	Basic Local Alignment Search Tool
326	cp	complete
327	cpDNA	chloroplast DNA
328	IR	Inverted Repeat
329	LSC	Large single-copy region
330	MAP	Maximum a posteriori
331	ML	Maximum Likelihood
332	MP	Maximum Parsimony
333	NGS	Next Generation Sequencing
334	NJ	Neighbor-joining
335	OTU	Operational Taxonomic Units
336	PCR	Polymerase Chain Reaction
337	RADP	Random Amplified Polymorphic DNA
338	RFLP	Restriction Fragment Length Polymorphism
339	SCAR	Sequence Characterized Amplified Region
340	SNP	Single Nucleotide Polymorphism
341	SSC	Small single-copy region
342	SSR	Simple Sequence Repeat
343	UPGMA	Unweighted Pair Group Method with Arithmetic mean

344

345

346

347

348

349

350

351 **REFERENCES**

- 352 1. Kress WJ, Wurdack KJ, Zimmer EA, Weigt LA, Janzen DH. Use of DNA barcodes to identify
353 flowering plants. Proceedings of the National Academy of Sciences of the United States of
354 America. 2005;102(23):8369-8374.

- 355 2. Kim HM, Oh SH, Bhandari GS, Kim CS, Park CW. DNA barcoding of Orchidaceae in Korea.
356 Molecular Ecology Resources. 2014;14(3):499-507.
- 357 3. Abbas B, Dailami M, Listyorini FH, Munarti. Genetic Variations and Relationships of Papua's
358 Endemic Orchids Based on RAPD Markers. Natural Science. 2017;9:377-385.
- 359 4. Kim JS, Kim HT, Son S-W, Kim J-H. Molecular identification of endangered Korean lady's
360 slipper orchids (*Cypripedium*, Orchidaceae) and related taxa. Botany. 2015;93(9):603-610.
- 361 5. Peyachoknagul S, Mongkolsiriwatana C, Wannapinpong S, Huehne PS, Srikulnath K.
362 Identification of native *Dendrobium* species in Thailand by PCR-RFLP of rDNA-ITS and
363 chloroplast DNA. ScienceAsia. 2014;40:113-120.
- 364 6. Sun YW, Liao YJ, Hung YS, Chang JC, Sung JM. Development of ITS sequence based SCAR
365 markers for discrimination of *Paphiopedilum armeniacum*, *Paphiopedilum micranthum*,
366 *Paphiopedilum delenatii* and their hybrids. Scientia Horticulturae. 2011;127(3):405-410.
- 367 7. Pereira F, Carneiro J, Amorim A. Identification of species with DNA-based technology:
368 current progress and challenges. Recent patents on DNA and gene sequences.
369 2008;2(3):187-99.
- 370 8. Yang F, Ding F, Chen H, He M, Zhu S, Ma X, Jiang L, Li H. DNA Barcoding for the
371 Identification and Authentication of Animal Species in Traditional Medicine. Evidence-
372 Based Complementary and Alternative Medicine. 2018;2018:18.
- 373 9. Waugh J. DNA barcoding in animal species: progress, potential and pitfalls. Bioessays.
374 2007;29(2):188-97.
- 375 10. Li X, Wu L, Wang J, Sun J, Xia X, Geng X, Wang X, Xu Z, Xu Q. Genome sequencing of rice
376 subspecies and genetic analysis of recombinant lines reveals regional yield- and quality-
377 associated loci. BMC Biology. 2018;16(1):102.
- 378 11. Ray S and Satya P. Next generation sequencing technologies for next generation plant
379 breeding. Frontiers in plant science. 2014;5:367-367.
- 380 12. Thottathil GP, Jayasekaran K, Othman AS. Sequencing Crop Genomes: A Gateway to
381 Improve Tropical Agriculture. Tropical life sciences research. 2016;27(1):93-114.
- 382 13. Zhou X, Bai X, Xing Y. A Rice Genetic Improvement Boom by Next Generation Sequencing.
383 Current Issues in Molecular Biology. 2018;27:109-126.
- 384 14. Hebert PD, Cywinska A, Ball SL, deWaard JR. Biological identifications through DNA
385 barcodes. Proceedings of the Royal Society. B, Biological sciences. 2003;270(1512):313-21.
- 386 15. Das S and Deb B. DNA barcoding of fungi using Ribosomal ITS Marker for genetic diversity
387 analysis: A Review. International Journal of Pure and Applied Bioscience. 2015;3(3):160-
388 167.
- 389 16. Selvaraj D, Park JI, Chung MY, Cho YG, Ramalingam S, Nou I-S, Utility of DNA Barcoding for
390 Plant Biodiversity Conservation. Vol. 1. 2013.
- 391 17. Techen N, Parveen I, Pan Z, Khan IA. DNA barcoding of medicinal plant material for
392 identification. Current Opinion in Biotechnology. 2014;25:103-110.
- 393 18. Teixeira da Silva JA, Jin X, Dobranszki J, Lu J, Wang H, Zotz G, Cardoso JC, Zeng S. Advances
394 in *Dendrobium* molecular research: Applications in genetic variation, identification and
395 breeding. Molecular Phylogenetics and Evolution. 2016;95:196-216.
- 396 19. Vu H-T, Huynh P, Tran H-D, Le L. In Silico Study on Molecular Sequences for Identification
397 of *Paphiopedilum* Species. Evolutionary Bioinformatics. 2018;14:117693431877454.
- 398 20. Fiser Pecnikar Z and Buzan EV. 20 years since the introduction of DNA barcoding: from
399 theory to application. Journal of Applied Genetics. 2014;55(1):43-52.
- 400 21. Hollingsworth PM. Refining the DNA barcode for land plants. Proceedings of the National
401 Academy of Sciences of the United States of America. 2011;108(49):19451-2.
- 402 22. Saddhe AA and Kumar K. DNA barcoding of plants: Selection of core markers for taxonomic
403 groups. Plant Science Today; Vol 5 No 1 (2018). 2017.
- 404 23. Duan H, Chen F, Liu W, Zhou C, Zhou Y. Research and Applications of DNA Barcode in
405 Identification of Plant Species. Research in Plant Biology. 2014;4(3).
- 406 24. Ganie SH, Upadhyay P, Das S, Prasad Sharma M. Authentication of medicinal plants by DNA
407 markers. Plant Gene. 2015;4:83-99.

- 408 25. Parveen I, Singh HK, Malik S, Raghuvanshi S, Babbar SB. Evaluating five different loci (rbcl,
409 rpoB, rpoC1, matK, and ITS) for DNA barcoding of Indian orchids. *Genome*. 2017;60(8):665-
410 671.
- 411 26. Rajaram MC, Yong C, Azlan GJ, Go R. DNA Barcoding of Endangered *Paphiopedilum* species
412 (Orchidaceae) of Peninsular Malaysia. *Phytotaxa*. 2019;387(2):94-104.
- 413 27. Chattopadhyay P, Banerjee G, Banerjee N. Distinguishing Orchid Species by DNA
414 Barcoding: Increasing the Resolution of Population Studies in Plant Biology. *Omic*.
415 2017;21(12):711-720.
- 416 28. Feng S, Jiang Y, Wang S, Jiang M, Chen Z, Ying Q, Wang H. Molecular Identification of
417 *Dendrobium* Species (Orchidaceae) Based on the DNA Barcode ITS2 Region and Its
418 Application for Phylogenetic Study. *International journal of molecular sciences*.
419 2015;16(9):21975-21988.
- 420 29. Wu C-T, Gupta SK, Wang AZ-M, Lo S-F, Kuo C-L, Ko Y-j, Chen C-L, Hsieh C-C, Tsay H-S.
421 Internal Transcribed Spacer Sequence Based Identification and Phylogenic Relationship of
422 Herba Dendrobii. *Journal of Food and Drug Analysis*. 2012;20(1):143-151.
- 423 30. Ginibun FC, Saad MRM, Hong TL, Othman RY, Khalid N, Bhasu S. Chloroplast DNA
424 Barcoding of *Spathoglottis* Species for Genetic Conservation. *Acta Horticulturae*.
425 2010;878:453-460.
- 426 31. Singh HK, Parveen I, Raghuvanshi S, Babbar SB. The loci recommended as universal
427 barcodes for plants on the basis of floristic studies may not work with congeneric species
428 as exemplified by DNA barcoding of *Dendrobium* species. *BMC Research Notes*. 2012;5:42.
- 429 32. Tanee T, Chadmuk P, Sudmoon R, Chaveerach A, Noikotr K. Genetic analysis for
430 identification, genomic template stability in hybrids and barcodes of the *Vanda* species
431 (Orchidaceae) of Thailand. *African Journal of Biotechnology*. 2012;11(55):11772-11781.
- 432 33. Yao H, Song JY, Ma XY, Liu C, Li Y, Xu HX, Han JP, Duan LS, Chen SL. Identification of
433 *Dendrobium* species by a candidate DNA barcode sequence: the chloroplast psbA-trnH
434 intergenic region. *Planta Medica*. 2009;75(6):667-9.
- 435 34. Givnish TJ, Spalink D, Ames M, Lyon SP, Hunter SJ, Zuluaga A, Iles WJD, Clements MA,
436 Arroyo MTK, Leebens-Mack J, Lorena E, Ricardo K, Kurt MN, Whitten WM, Norris HW,
437 Kenneth MC. Orchid phylogenomics and multiple drivers of their extraordinary
438 diversification. *Proceedings of the Royal Society B*. 2015;282:20151553.
- 439 35. Santos C and Pereira F. Identification of plant species using variable length chloroplast
440 DNA sequences. *Forensic Science International: Genetics*. 2018;36:1-12.
- 441 36. Asahina H, Shinozaki J, Masuda K, Morimitsu Y, Satake M. Identification of medicinal
442 *Dendrobium* species by phylogenetic analyses using matK and rbcl sequences. *Journal of*
443 *Natural Medicines*. 2010;64(2):133-8.
- 444 37. Tang Y, Yukawa T, Bateman RM, Jiang H, Peng H. Phylogeny and classification of the East
445 Asian *Amitostigma alliance* (Orchidaceae: Orchideae) based on six DNA markers. *BMC*
446 *Evolutionary Biology*. 2015;15:96.
- 447 38. Tran DD, Khuat HT, La TN, Nguyen TTT, Pham BH, Nguyen TK, Tran HD, Do MT, Tran DK.
448 Identification of Vietnamese Native *Dendrobium* Species Based on Ribosomal DNA Internal
449 Transcribed Spacer Sequence. *Advanced Studies in Biology*. 2018;10(1):1-12.
- 450 39. Poovitha S, Stalin N, Balaji R, Parani M. Multi-locus DNA barcoding identifies matK as a
451 suitable marker for species identification in *Hibiscus* L. *Genome*. 2016;59(12):1150-1156.
- 452 40. EMBL-EBI. FASTA Help and Documentation. 2019. Accessed
453 <https://www.ebi.ac.uk/seqdb/confluence/display/JDSAT/FASTA+Help+and+Documentatio>
454 [n](https://www.ebi.ac.uk/seqdb/confluence/display/JDSAT/FASTA+Help+and+Documentation).
- 455 41. Lloyd A. The Clustal W WW server at the EBI. 1997. Accessed
456 http://www.ebi.ac.uk/embnet.news/vol4_3/clustalw1.html.
- 457 42. Katoh K, Asimenos G, Toh H. Multiple alignment of DNA sequences with MAFFT. *Methods*
458 *in molecular biology* (Clifton, N.J.). 2009;537:39-64.
- 459 43. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
460 *Nucleic Acids Res*. 2004;32(5):1792-7.

- 461 44. Katoh K, Rozewicki J, Yamada KD. MAFFT online service: multiple sequence alignment,
462 interactive sequence choice and visualization. *Briefings in Bioinformatics*. 2017;20(4):1160-
463 1166.
- 464 45. Ghorbani A, Gravendeel B, Selliah S, Zarre S, de Boer H. DNA barcoding of tuberous
465 Orchidoideae: a resource for identification of orchids used in Salep. *Molecular Ecology*
466 *Resources*. 2017;17(2):342-352.
- 467 46. Guo YY, Huang LQ, Liu ZJ, Wang XQ. Promise and Challenge of DNA Barcoding in Venus
468 Slipper (*Paphiopedilum*). *PLOS ONE*. 2016;11(1):e0146880.
- 469 47. Meier R, Shiyang K, Vaidya G, Ng PK. DNA barcoding and taxonomy in Diptera: a tale of
470 high intraspecific variability and low identification success. *Systematic Biology*.
471 2006;55(5):715-28.
- 472 48. Xiang XG, Hu H, Wang W, Jin XH. DNA barcoding of the recently evolved genus
473 *Holcoglossum* (Orchidaceae: Aeridinae): a test of DNA barcode candidates. *Molecular*
474 *Ecology Resources*. 2011;11(6):1012-21.
- 475 49. Xu S, Li D, Li J, Xiang X, Jin W, Huang W, Jin X, Huang L. Evaluation of the DNA Barcodes in
476 *Dendrobium* (Orchidaceae) from Mainland Asia. *PLOS ONE*. 2015;10(1):e0115168.
- 477 50. Gao T, Yao H, Song J, Liu C, Zhu Y, Ma X, Pang X, Xu H, Chen S. Identification of medicinal
478 plants in the family Fabaceae using a potential DNA barcode ITS2. *Journal of*
479 *Ethnopharmacology*. 2010;130(1):116-21.
- 480 51. Siripiyasing P. DNA barcoding of the *Cymbidium* species (Orchidaceae) in Thailand. *African*
481 *journal of agricultural research*. 2012;7(3):393-404.
- 482 52. Meier R, Zhang G, Ali F. The Use of Mean Instead of Smallest Interspecific Distances
483 Exaggerates the Size of the “Barcoding Gap” and Leads to Misidentification. *Systematic*
484 *Biology*. 2008;57(5):809-813.
- 485 53. Guo Y-Y, Luo Y-B, Liu Z-J, Wang X-Q. Evolution and Biogeography of the Slipper Orchids:
486 Eocene Vicariance of the Conduplicate Genera in the Old and New World Tropics. *PLOS*
487 *ONE*. 2012;7(6):e38788.
- 488 54. Shaw J, Lickey EB, Schilling EE, Small RL. Comparison of whole chloroplast genome
489 sequences to choose noncoding regions for phylogenetic studies in angiosperms: the
490 tortoise and the hare III. *American Journal of Botany*. 2007;94(3):275-88.
- 491 55. Trung KH, Khanh TD, Ham LH, Duong TD, Khoa T. Molecular Phylogeny of the Endangered
492 Vietnamese *Paphiopedilum* Species Based on the Internal Transcribed Spacer of the
493 Nuclear Ribosomal DNA. *Advanced Studies in Biology*. 2013;5(7):337 - 346.
- 494 56. Meyer CP and Paulay G. DNA barcoding: error rates based on comprehensive sampling.
495 *PLoS biology*. 2005;3(12):e422-e422.
- 496 57. Gigot G, Van Alphen-Stahl J, Bogarin D, Warner J, Chase MW, Savolainen V. Finding a
497 suitable DNA barcode for Mesoamerican orchids. *Lankesteriana International Journal on*
498 *Orchidology*. 2007;7(1-2):200-203.
- 499 58. Parveen I, Singh HK, Raghuvanshi S, Pradhan UC, Babbar SB. DNA barcoding of endangered
500 Indian *Paphiopedilum* species. *Molecular Ecology Resources*. 2012;12(1):82-90.
- 501 59. Yukawa T, Kinoshita1 A, Tanaka N. Molecular Identification Resolves Taxonomic Confusion
502 in *Grammatophyllum speciosum* Complex (Orchidaceae). *Bulletin of the National Museum*
503 *of Nature and Science, Series B*. 2013;39(3):137–145.
- 504 60. Chen L-J, Liu Z-J, Li Y-Y, Li L-Q. A new orchid *Paphiopedilum guangdongense* and its
505 molecular evidence. *Journal of Systematics and Evolution*. 2010;48(5):350-355.
- 506 61. Lahaye R, van der Bank M, Bogarin D, Warner J, Pupulin F, Gigot G, Maurin O, Duthoit S,
507 Barraclough TG, Savolainen V. DNA barcoding the floras of biodiversity hotspots.
508 *Proceedings of the National Academy of Sciences of the United States of America*.
509 2008;105(8):2923-8.
- 510 62. Slabbinck B, Dawyndt P, Martens M, De Vos P, De Baets B. TaxonGap: a visualization tool
511 for intra- and inter-species variation among individual biomarkers. *Bioinformatics*.
512 2008;24(6):866-7.

- 513 63. Nei M and Saitou N. The neighbor-joining method: a new method for reconstructing
514 phylogenetic trees. *Molecular Biology and Evolution*. 1987;4(4):406-425.
- 515 64. Felsenstein J. Evolutionary trees from DNA sequences: A maximum likelihood approach.
516 *Journal of Molecular Evolution*. 1981;17(6):368-76.
- 517 65. Penny D. Inferring Phylogenies.—Joseph Felsenstein. 2003. Sinauer Associates,
518 Sunderland, Massachusetts. *Systematic Biology*. 2004;53(4):669-670.
- 519 66. Felsenstein J. Confidence limits on phylogenies: An approach using the bootstrap.
520 *Evolution*. 1985;39(4):783-791.
- 521 67. Mar JC, Harlow TJ, Ragan MA. Bayesian and maximum likelihood phylogenetic analyses of
522 protein sequence data under relative branch-length differences and model violation. *BMC*
523 *evolutionary biology*. 2005;5:8-8.
- 524 68. Kolaczowski B and Thornton JW. Performance of maximum parsimony and likelihood
525 phylogenetics when evolution is heterogeneous. *Nature*. 2004;431:980.
- 526 69. Kumar S, Nei M, Dudley J, Tamura K. MEGA: a biologist-centric software for evolutionary
527 analysis of DNA and protein sequences. *Briefings in Bioinformatics*. 2008;9(4):299-306.
- 528 70. Swofford DL. PAUP*. Phylogenetic analysis using parsimony and other methods. Version
529 4.0. 2003.
- 530 71. Posada D. jModelTest: phylogenetic model averaging. *Mol Biol Evol*. 2008;25(7):1253-6.
- 531 72. Huelsenbeck JP and Ronquist F. MRBAYES: Bayesian inference of phylogenetic trees.
532 *Bioinformatics*. 2001;17(8):754-5.
- 533 73. Xu Q, Zhang G-Q, Liu Z-J, Luo Y-B. Two new species of *Dendrobium* (Orchidaceae:
534 Epidendroideae) from China: evidence from morphology and DNA. *Phytotaxa*.
535 2014;174(3):15.
- 536 74. Xiang X-G, Jin W-T, Li D-Z, Schuiteman A, Huang W-C, Li J-W, Jin X-H, Li Z-Y. Phylogenetics
537 of Tribe Collabieae (Orchidaceae, Epidendroideae) Based on Four Chloroplast Genes with
538 Morphological Appraisal. *PLOS ONE*. 2014;9(1):e87625.
- 539 75. Yang J-B, Yang S-X, Li H-T, Yang J, Li D-Z. Comparative chloroplast genomes of *Camellia*
540 species. *PloS one*. 2013;8(8):e73053-e73053.
- 541 76. Nock CJ, Waters DL, Edwards MA, Bowen SG, Rice N, Cordeiro GM, Henry RJ. Chloroplast
542 genome sequences from total DNA for plant identification. *Plant Biotechnology Journal*.
543 2011;9(3):328-33.
- 544 77. Kane N, Sveinsson S, Dempewolf H, Yang JY, Zhang D, Engels JM, Cronk Q. Ultra-barcoding
545 in cacao (*Theobroma* spp.; Malvaceae) using whole chloroplast genomes and nuclear
546 ribosomal DNA. *American Journal of Botany*. 2012;99(2):320-9.
- 547 78. Parker J, Helmstetter AJ, Devey D, Wilkinson T, Papadopulos AST. Field-based species
548 identification of closely-related plants using real-time nanopore sequencing. *Scientific*
549 *Reports*. 2017;7(1):8345.
- 550 79. Day PD, Berger M, Hill L, Fay MF, Leitch AR, Leitch IJ, Kelly LJ. Evolutionary relationships in
551 the medicinally important genus *Fritillaria* L. (Liliaceae). *Molecular Phylogenetics and*
552 *Evolution*. 2014;80:11-9.
- 553 80. Turktas M, Aslay M, Kaya E, Ertugrul F. Molecular characterization of phylogenetic
554 relationships in *Fritillaria* species inferred from chloroplast trnL-trnF sequences. *Turkish*
555 *Journal of Biology*. 2012;36(5):552-560.
- 556 81. Khourang M, Babaei A, Sefidkon F, Naghavi MR, Asgari D, Potter D. Phylogenetic
557 relationship in *Fritillaria* spp. of Iran inferred from ribosomal ITS and chloroplast trnL-trnF
558 sequence data. *Biochemical Systematics and Ecology*. 2014;57:451-457.
- 559 82. Bi Y, Zhang MF, Xue J, Dong R, Du YP, Zhang XH. Chloroplast genomic resources for
560 phylogeny and DNA barcoding: a case study on *Fritillaria*. *Scientific Reports*.
561 2018;8(1):1184.
- 562 83. Chen X, Zhou J, Cui Y, Wang Y, Duan B, Yao H. Identification of *Ligularia* Herbs Using the
563 Complete Chloroplast Genome as a Super-Barcode. *Frontiers in pharmacology*.
564 2018;9:695-695.

- 565 84. Lee J, Chon J, Lim J, Kim E-K, Nah G. Characterization of Complete Chloroplast Genome of
566 *Allium victorialis* and Its Application for Barcode Markers. *Plant Breeding and*
567 *Biotechnology*. 2017;5(3):221-227.
- 568 85. Dong W, Liu H, Xu C, Zuo Y, Chen Z, Zhou S. A chloroplast genomic strategy for designing
569 taxon specific DNA mini-barcodes: a case study on ginsengs. *BMC genetics*. 2014;15:138-
570 138.
- 571 86. Zhou Y, Nie J, Xiao L, Hu Z, Wang B. Comparative Chloroplast Genome Analysis of Rhubarb
572 Botanical Origins and the Development of Specific Identification Markers. *Molecules*.
573 2018;23(11).
- 574 87. Curci PL, De Paola D, Danzi D, Vendramin GG, Sonnante G. Complete chloroplast genome
575 of the multifunctional crop globe artichoke and comparison with other Asteraceae. *PLoS*
576 *One*. 2015;10(3):e0120589.
- 577 88. Yi D-K, Choi K, Joo M, Yang JC, Mustafina FU, Han J-S, Son DC, Chang KS, Shin CH, Lee Y-M.
578 The complete chloroplast genome sequence of *Abies nephrolepis* (Pinaceae: Abietoideae).
579 *Journal of Asia-Pacific Biodiversity*. 2016;9(2):245-249.
- 580 89. Zhang Y, Guan W, Zhang X, Li L. The Complete Chloroplast Genomes of Asteraceae Species.
581 *Research & Reviews: Journal of Botanical Sciences*. 2016;5(1):24-28.
- 582 90. Jheng CF, Chen TC, Lin JY, Chen TC, Wu WL, Chang CC. The comparative chloroplast
583 genomic analysis of photosynthetic orchids and developing DNA markers to distinguish
584 *Phalaenopsis* orchids. *Plant Science*. 2012;190:62-73.
- 585 91. Lin J-Y, Lin B-Y, Chang C-D, Liao S-C, Liu Y-C, Wu W-L, Chang C-C. Evaluation of chloroplast
586 DNA markers for distinguishing *Phalaenopsis* species. *Scientia Horticulturae*.
587 2015;192:302-310.
- 588 92. Yu XQ, Drew BT, Yang JB, Gao LM, Li DZ. Comparative chloroplast genomes of eleven
589 *Schima* (Theaceae) species: Insights into DNA barcoding and phylogeny. *PLOS ONE*.
590 2017;12(6):e0178026.
- 591 93. Peyachoknagul S, Mongkolsiriwatana C, Wannapinpong S, Huehne P, Srikulnath K.
592 Identification of native *Dendrobium* species in Thailand by PCR-RFLP of rDNA-ITS and
593 chloroplast DNA. *ScienceAsia*. 2014;40(2):113–120.
- 594 94. Boer HJ, Ghorbani A, Manzanilla V, Raclariu AC, Kreziou A, Ounjai S, Osathanunkul M,
595 Gravendeel B. DNA metabarcoding of orchid-derived products reveals widespread illegal
596 orchid trade. *Proceedings of the Royal Society B: Biological Sciences*. 2017;284(1863).
- 597 95. Veldman S, Kim SJ, van Andel TR, Bello Font M, Bone RE, Bytebier B, Chuba D, Gravendeel
598 B, Martos F, Mpatwa G, Ngugi G, Vinya R, Wightman N, Yokoya K, de Boer HJ. Trade in
599 Zambian Edible Orchids-DNA Barcoding Reveals the Use of Unexpected Orchid Taxa for
600 Chikanda. *Genes (Basel)*. 2018;9(12).

601