# Feasibility of nuclear ribosomal region ITS1 over ITS2 in barcoding taxonomically challenging genera of subtribe Cassiinae (Fabaceae)

Priyanka Mishra [1], Amit Kumar [1], Vereena Rodrigues [1], Ashutosh K Shukla [2], Velusamy Sundaesan [Corresp. 1]

[1] Department of Plant Biology & Systematics, Central Institute of Medicinal and Aromatics Plant Research Center, Bangalore, Karnataka, India

[2] Biotechnology Division, CSIR-Central Institute of Medicinal and Aromatic Plants, Lucknow, Uttar Pradesh, India

Corresponding Author: Velusamy Sundaesan
Email address: vsundaresan@cimap.res.in

**Premise of the Study.** The internal transcribed spacer (ITS) region is situated between 18S and 26S in a polycistronic rRNA precursor transcript. It had been proved to be the most commonly sequenced region across plant species to resolve phylogenetic relationships ranging from shallow to deep taxonomic levels. Despite several taxonomical revisions in Cassiinae, a stable phylogeny remains elusive at the molecular level, particularly concerning the delineation of species in the genera *Cassia, Senna* and *Chamaecrista*. This study addresses the comparative potential of ITS datasets (ITS1, ITS2 and concatenated) in resolving the underlying morphological disparity in the highly complex genera, to assess their discriminatory power as potential barcode candidates in Cassiinae.

**Methodology.** A combination of experimental data and an in-silico approach based on threshold genetic distances, sequence similarity based and hierarchical tree-based methods was performed to decipher the discriminating power of ITS datasets on 18 different species of Cassiinae complex**.** Lab-generated **s**equences were compared against those available in the GenBank using BLAST and were aligned through MUSCLE 3.8.31 and analysed in PAUP 4.0 and BEAST1.8 using parsimony ratchet, maximum likelihood and Bayesian inference (BI) methods of gene and species tree reconciliation with bootstrapping. DNA barcoding gap was realized based on the Kimura two-parameter distance model (K2P) in TaxonDNA and MEGA.

**Principal Findings.** Based on the K2P distance, significant divergences between the inter- and intra-specific genetic distances were observed, while the presence of a DNA barcoding gap was obvious. The ITS1 region efficiently identified 81.63% and 90% of species using TaxonDNA and BI methods, respectively. The PWG-distance method based on simple pairwise matching indicated the significance of ITS1 whereby highest number of variable (210) and informative sites (206) were obtained. The BI tree-based methods outperformed the similarity-based methods producing well-resolved phylogenetic trees with many nodes well supported by bootstrap analyses.

**Conclusion.** The reticulated phylogenetic hypothesis using the ITS1 region mainly supported the relationship between the species of Cassiinae established by traditional morphological methods. The ITS1 region showed a higher discrimination power and desirable characteristics as compared to ITS2 and ITS1+2 there by concluding to be the locus of choice. Considering the complexity of the group and the underlying biological ambiguities, the results presented here are encouraging for developing DNA barcoding as a useful tool for resolving taxonomical challenges in corroboration with morphological framework.

1 **Feasibility of nuclear ribosomal region ITS1 over ITS2 in barcoding**

2 **taxonomically challenging genera of subtribe Cassiinae (Fabaceae)**

3 Priyanka Mishra[1], Amit Kumar[1], Vereena Rodrigues[1], Ashutosh K. Shukla[2] and Velusamy

4 Sundaresan[1,*]

5 [1] Department of Plant Biology & Systematics, CSIR-Central Institute of Medicinal and Aromatic

6 Plants, Research Center, Bangalore – 560065, Karnataka, India

7 [2] Biotechnology Division, CSIR - Central Institute of Medicinal and Aromatic Plants, Lucknow –

8 226015, Uttar Pradesh, India

9 [*]**Corresponding author:**

10 Velusamy Sundaresan, Ph.D.

11 E-mail: vsundaresan@cimap.res.in, resanvs@gmail.com

13 **ABSTRACT**

14 **Premise of the Study.** The internal transcribed spacer (ITS) region is situated between 18S and

15 26S in a polycistronic rRNA precursor transcript. It had been proved to be the most commonly

16 sequenced region across plant species to resolve phylogenetic relationships ranging from shallow

17 to deep taxonomic levels. Despite several taxonomical revisions in Cassiinae, a stable phylogeny

18 remains elusive at the molecular level, particularly concerning the delineation of species in the

19 genera *Cassia, Senna* and *Chamaecrista*. This study addresses the comparative potential of ITS

20 datasets (ITS1, ITS2 and concatenated) in resolving the underlying morphological disparity in

21 the highly complex genera, to assess their discriminatory power as potential barcode candidates

22 in Cassiinae.

23 **Methodology.** A combination of experimental data and an in-silico approach based on threshold

24 genetic distances, sequence similarity based and hierarchical tree-based methods was performed

25 to decipher the discriminating power of ITS datasets on 18 different species of Cassiinae

26 complex**.** Lab-generated **s**equences were compared against those available in the GenBank using

27 BLAST and were aligned through MUSCLE 3.8.31 and analysed in PAUP 4.0 and BEAST1.8

28 using parsimony ratchet, maximum likelihood and Bayesian inference (BI) methods of gene and

29 species tree reconciliation with bootstrapping. DNA barcoding gap was realized based on the

30 Kimura two-parameter distance model (K2P) in TaxonDNA and MEGA.

31 **Principal Findings.** Based on the K2P distance, significant divergences between the inter- and

32 intra-specific genetic distances were observed, while the presence of a DNA barcoding gap was

33 obvious. The ITS1 region efficiently identified 81.63% and 90% of species using TaxonDNA

34 and BI methods, respectively. The PWG-distance method based on simple pairwise matching

35 indicated the significance of ITS1 whereby highest number of variable (210) and informative

36    sites (206) were obtained. The BI tree-based methods outperformed the similarity-based methods

37    producing well-resolved phylogenetic trees with many nodes well supported by bootstrap

38    analyses.

39    **Conclusion.** The reticulated phylogenetic hypothesis using the ITS1 region mainly supported the

40    relationship between the species of Cassiinae established by traditional morphological methods.

41    The ITS1 region showed a higher discrimination power and desirable characteristics as compared

42    to ITS2 and ITS1+2 there by concluding to be the locus of choice. Considering the complexity of

43    the group and the underlying biological ambiguities, the results presented here are encouraging

44    for developing DNA barcoding as a useful tool for resolving taxonomical challenges in

45    corroboration with morphological framework.

## INTRODUCTION

DNA barcoding is an important tool for research in biodiversity hot-spots based on the identification and standardization of specific region of the plant genome that can be sequenced routinely in diverse sample sets to identify and discriminate species from one another (*Hebert et al., 2003; Gregory, 2005*). The revolution introduced by DNA barcoding relies on molecularization (variability in molecular markers), computerization (transposition of the data through bioinformatics workbench) and standardization (extension of approach to diverse group) of traditional taxonomical framework to easily associate all life stages of a biological entity (*Casiraghi et al., 2010*). The short, variable and standardized DNA sequence can be termed as DNA barcode when it mirrors the distributions of intra- and inter-specific variabilities separated by a distance called 'DNA barcoding gap' and characterizes conserved flanking regions for development of universal primers across highly divergent taxa (*Kress et al., 2005; Savolainen et al., 2005; Hollingsworth et al. 2009*).

In the past, DNA barcoding in plants has been extensively reviewed (*Vijayan & Tsou, 2010; Hollingsworth et al., 2011*), but still there is a considerable debate on the consensus of the choice of a standard region (*Mishra et al., 2015*). Apart from the accepted mitochondrial cytochrome oxidase I gene (*COI*) in animals and the nuclear ribosomal internal transcribed spacer (ITS) region in fungi, the search for an analogous region in plants focused attention on the plastid genome (*Chase et al., 2005; Kress et al., 2005; Nilsson et al., 2006; Fazekas et al., 2009*). Subsequently, major individual candidate regions *matK, rbcL, rpoB, rpoC1*, and the intergenic spacers ITS, *trnH-psbA*, *trnL-F*, *atpF-atpH* and *psbK-psbI*, etc. were tested for use in plants on their discrimination capacity. Due to pitfalls and challenges associated with a single locus, the combination of loci emerged as a promising choice to obtain appropriate species

70    discrimination (*Chase et al., 2007; Kress & Erickson, 2007; Fazekas et al., 2008; CBOL Plant*

71    *Working Group, 2009; Hollingsworth et al., 2011*).

72          The ITS region in plants has been shown to perform as a powerful phylogenetic marker

73    when compared with either coding or noncoding plastid markers due to high copy number of

74    rRNA genes and high degree of variations even between the closely related species (*Álvarez &*

75    *Wendel, 2003; Chase et al., 2007; China Plant BOL Group, 2011; Li et al., 2014*). The

76    availability of several universal primer sets and moderate size of 500–750 bp provides an

77    advantageous feature in deciphering the riddles within and among various taxa. The spacer DNA

78    occurs as intercalated in the 16S–5.8S–26S region of rDNA locus and consists of ITS1, ITS2 and

79    the highly conserved 5.8S. Also, many studies have compared the discriminatory power of ITS

80    region in its entirety with ITS2, proposing ITS2 as an alternative barcode to entire ITS region

81    because of sufficient variation in primary sequences and secondary structures (*Chen et al., 2010;*

82    *Gao et al., 2010; Han et al., 2013*). Despite the problems in amplifying and directly sequencing

83    the entire region, ITS1 has been tested as a better barcode for eukaryotic species (*Wang et al.,*

84    *2014*) and also a successful region for the members of legume family (*Yadav et al., 2016*).

85          Fabaceae (Legumes) are the third largest family of flowering plants with Caesalpinioidae

86    being the second largest of the three subfamilies (*Irwin & Barneby, 1981*). Cassiinae is a subtribe

87    of Fabaceae in the subfamily Caesalpinioidae, comprising of three genera, viz. *Cassia* L. sens.

88    str., *Senna* P. Mill., and *Chamaecrista* Moench. Genus *Cassia* L. sens. *lat.*, is one of the twenty-

89    five largest genera of dicotyledonous plant with high diversity of secondary metabolites which

90    serve as medicinal, nutraceuticals and sustainable agriculture etc. (*Singh, 2001*). Tinnevelly

91    *Senna* is the second largest exported herb drug in the country and contributes significantly in the

92    range of 5000 metric tons per year as commercial products (*Seethapathy et al., 2014*). Despite

93   several studies by many taxonomists, either on the whole family or at the genus level, there has

94   been considerable divergence of opinion concerning the delimitations and taxonomic status of

95   the subgenera at the molecular level. The wide variability in habit ranging from tall trees to

96   delicate annual herbs, floral and vegetative features, pods variability etc had made its

97   taxonomical framework quite complex and intriguing (*Singh, 2001*). Cytological and

98   karyological studies of 17 taxa of *Cassia*, showed no correlation between the habit and karyotype

99   symmetry of various species (*Bir & Kumari, 1982*). Thus the identification of the species has

100  proved tricky and is rather difficult to account for the entire genetic variation existing in the

101  genera. A robust and reliable method is crucial to discriminate plant species to secure their

102  diversity.

103       Few studies in *Cassia* have been conducted utilizing the dominant molecular markers

104  (*Mohanty et al., 2010*), plastid and nuclear region markers for different purposes

105  (*Purushothaman et al., 2014*; *Seethapathy et al., 2014*). The studies demonstrated the subsequent

106  contribution of markers in assessing product adulteration in herbal drug market in India

107  (*Seethapathy et al., 2014*). Although the results were not based on evolutionary relationships

108  concept, they did indicate a potential role of different regions (markers) in resolving species

109  complexity in *Cassia* (*Mohanty et al., 2010; Purushothaman et al., 2014*).

110       In this study, we evaluated the potential ability of ITS regions for identifying and

111  discriminating subtribe Cassiinae based on a representative sample consisting of approximately

112  half of the genera. The applicability and effectiveness of ITS regions (ITS1 and ITS2) in

113  discriminating species across the genera *Cassia, Senna* and *Chamaecrista* were studied for the

114  first time. The sufficient sequences available in GenBank with nuclear region ITS were included

115  for analysis. The main goals of this study were as follows: (i) to infer applicability and efficacy

116    of the ITS regions (ITS1, ITS2 and ITS1+2) as barcoding candidates for subtribe *Cassiinae*; (ii)

117    to test the reliability of the underlying taxonomic monographs at the genome level in resolving

118    congeneric species; and (iii) to compare different methods of evaluating DNA barcodes in these

119    highly complex genera.

120    **MATERIALS AND METHODS**

121    **Taxon sampling, DNA amplification and sequencing**

122    A total of 54 accessions of 18 species belonging to three genera viz. *Cassia, Senna,* and

123    *Chamaecrista* from India were examined during the study. For obtaining the sequences

124    generated from molecular experiments in our lab, a total of 18 individuals corresponding to three

125    different genera were collected from different geographical regions of South Western Ghats and

126    Uttar Pradesh. The species were identified and authenticated using the morphological characters

127    described in a monographic study on Cassiinae in India (Singh, 2001) by Dr. V. Sundaresan,

128    Scientist, Central Institute of Medicinal and Aromatic Plants, Research Centre (Bangalore). For

129    each of the species, herbarium specimens were prepared and deposited at the Herbaria of the

130    Central Institute of Medicinal and Aromatic Plants (CIMAP Communication No.:

131    CIMAP/PUB/2016/24), Lucknow.

132        Legumes family produce a high diversity of secondary metabolites, which causes extreme

133    difficulty in isolation of high-quality nucleic acids. Based on literature and commercial kits

134    available, we attempted modification of several previously reported methods to isolate high

135    quality DNA. Ultimately, total genomic DNA from individual accessions was extracted from the

136    leaf tissues (dried in silica-gel) using the modified cetyl trimethyl ammonium bromide (CTAB)

137    protocol with necessary major modifications (*Khanuja et al., 1999*) and supplementing it with

138    the Nucleospin Plant II Maxi prep kit using the manufacturer's protocol (MACHEREY-NAGEL,

139    Duren Germany). The concentration of β-mercaptoethanol and PVP (Polyvinylpyrrolidone) were

140    increased to 2% v/v and 4% w/v, respectively. An additional chloroform-isoamyl alcohol (96:4)

141    purification step was performed to remove proteins and potentially interfering secondary

142    metabolites. Isolated DNA was checked for its quality and quantity by electrophoresis on a 0.8%

143    agarose gel and spectrophotometric analysis (NanoDrop, ND-1000, USA). The nuclear internal

144    transcribed spacer (ITS1 and ITS2) regions of all the individuals were amplified according to

145    PCR reaction conditions (94°C, 5 min; [30 cycles: 94°C, 1 min; 50°C, 1 min; 72°C, 1.5 min];

146    72°C, 7 min) following guidelines from the CBOL plant-working group and sequenced using

147    universal primers ITS5a forward 5'-CCTTATCATTTAGAGGAAGGAG-3' and ITS4 reverse 5'-

148    TCCTCCGCTTATTGATATGC-3' (*Kress et al., 2005*). PCR amplifications for each primer set

149    were carried out in a 50 µl volume solution containing 1x Taq DNA polymerase buffer, 200 µM

150    each dNTPs (dATP:dTTP:dCTP:dGTP in 1:1:1:1 parts), 10 pmol of each primer (forward and

151    reverse), 1 unit of Taq DNA polymerase and ≈25-50 ng of template DNA. The PCR fragment

152    lengths were determined on a 2% agarose gel. The PCR products were purified with Nucleospin

153    PCR purification kit (MACHEREY-NAGEL, Duren, Germany) as per the manufacturer's

154    instructions. Presence of the specific product was confirmed by running the purified PCR

155    products on 2% agarose gel. All the purified PCR products were subjected to double-stranded

156    sequencing using the Applied Biosystems Prism Big Dye Terminator Cycle Sequencing Kit

157    (Applied Biosystems, Foster City, CA) on an ABI 3130 XL automated sequencer (Applied

158    Biosystems).

159         Apart from the lab-generated sequences, all the nucleotide sequences belonging to genera

160    *Cassia, Senna,* and *Chamaecrista* for the regions ITS1 and ITS2 were downloaded from the

161  NCBI based on the blast results. The sequences were filtered on the basis of length (less than 300

162  bp were omitted), lack of voucher specimens as well as verification (sequences categorised as

163  unverified in GenBank were omitted). An effort was made to include minimum five individuals

164  for each species, but due to unavailability of sequences for few species in the NCBI database and

165  difficulty in obtaining the species in the field, the representatives of each species were limited to

166  three. The GenBank accession numbers used in this study are listed in Table 1.

167  **Data analysis**

168  Electropherograms corresponding to raw sequences of individual accessions from both the

169  forward and reverse primers were assembled and edited using CodonCode Aligner v.3.0.1

170  (CodonCode Corporation). Sequences were clipped at the end to avoid the presence of variable

171  sites introduced by the sequencing artefacts. Due to its well-conserved nature, the 5.8S gene

172  region was removed from any sequence so that the ITS1 and ITS2 regions could be analyzed

173  separately and concatenated. The edited sequences were then aligned with MUSCLE 3.8.31 on

174  the EMBLEBI website (http://www.ebi.ac.uk) with default parameter and adjusted manually in

175  BioEdit v7.1.3.0 (*Hall, 1999*). All the variable sites were rechecked on the original trace files. To

176  evaluate the effectiveness of ITS1, ITS2 and their combination (ITS1+2) as barcodes in the

177  concerned genera, three widely used methods viz. distance-based (PWG-distance), similarity-

178  based and tree-based were applied.

179  **Genetic Distance-Based Method**

180  To evaluate the measure of effective barcode locus, DNA barcoding gap was calculated using

181  TaxonDNA software with a 'pairwise summary' function under K2P nucleotide substitution

182  model (*Meier et al., 2006*). The pairwise genetic distance were calculated at the observed levels

183   of intra- and inter-specific divergenc for each barcode. To test the accurate species assignments,

184   the distributions of the pairwise intra- and inter-specific distances with 0.005 distance intervals

185   were generated. The histogram of distances vs. abundance were plotted to estimate the presence

186   of any barcoding gaps. For the PWG-distance method, the genetic pairwise distance was

187   estimated by MEGA version 6 (*Tamura et al., 2013*) using the Kimura two-parameter distance

188   model (K2P) with pairwise deletion of missing sites (*Kimura, 1980*). Average inter-specific

189   distance was used to characterize inter-specific divergence (*Meyer & Paulay, 2005, Meier et al.,*

190   *2008*) and 'all' intra-specific distance, mean 'theta' and coalescent depth were used to

191   characterize intra-specific distances. Finally, the obtained inter- and intra-specific distances were

192   plotted with frequency distribution in bin interval of 0.05 to illustrate the existing DNA

193   barcoding gap (*Meyer & Paulay, 2005, Lahaye et al., 2008*).

194   **DNA Sequence Similarity-Based Method**

195   To test the potentiality of ITS regions to identify species accurately based on sequence similarity,

196   the proportion of correct identifications were calculated using SpeciesIdentifier program from

197   the TAXONDNA software package with 'Best match' (BM), 'Best close match' (BCM) and 'All

198   species barcodes' functions (*Meyer & Paulay, 2005*). The tool examines all the sequences

199   present in aligned data set and compares each successive sequence with all the other sequences

200   to determine the closest match. The 'Best match' modules than classifies the sequences as

201   correct and incorrect based on the indicated pair from the similar species or different species

202   respectively. While the various equally best matches from different species are referred to be as

203   ambiguous. The 'Best close match' module works on the intra-species variability criterion and

204   considered to be the more rigorous method in TaxonDNA. The sequences classified as 'no

205   match' are the results above the calculated threshold value (*Meier et al., 2006*).

## Tree-Based Method

206

207 To evaluate the ability of candidate barcode to delimit the species into discrete clades or

208 monophyletic groups, three different optimality criteria (tree-building method) viz Neighbour-

209 joining with minimum evolution (NJ), maximum likelihood (ML) and Bayesian inference (BI)

210 were employed. To test the reliability of the result, NJ and ML trees were constructed and

211 compared with two different softwares: (i) In MEGA using the K2P distance as model of

212 substitution (*Tamura et al., 2013*) and (ii) In PAUP 4.0 with the HKY-gamma substitution model

213 (*Swofford, 2003*). The reliability of the node was assessed by a bootstrap test with 1000 pseudo-

214 replicates with the K2P distance options (*Felsenstein, 1988*). Bayesian sampling was performed

215 in BEAST1.8 using the operators: HKY substitution model with four gamma categories, a

216 constant-rate Yule tree prior and 10000 chain lengths and all other priors and operators with the

217 default settings. Coalescent tree priors were used for population-level analysis and speciation

218 prior were applied to estimate relationships and divergence times of inter-species data. Trees

219 were sampled for every 5000 generations resulting in a total of 10000 trees, and a burn-in of

220 5000000. Beast file was created using the BEAUti program v1.8.2 within Beast and performance

221 of each run was further analysed with the program Tracer (*Rambaut et al., 2012*). The resulting

222 Beast tree files were annotated through TreeAnnotator v1.8.2 and visualized and edited with

223 FigTree v1.4.2. (*Drummond et al., 2012, http://tree.bio.ed.ac.uk/software/figtree*). Visualization

224 and analysis of all the resulting trees through PAUP 4.0 was done in Dendroscope3 (*Huson &*

225 *Scornavacca, 2012*). Gaps were treated as missing data for all the phylogenetic analysis.

## RESULTS

226

## PCR amplification and sequence characteristics

227

228  The sequence characteristics of ITS regions evaluated in this study showed good success rates

229  (90%) for PCR amplification (ranging from 571bp - 1153bp with mean size ≈ 707bp ; gel images

230  can bé provided on request) and sequencing in both the direction using a single primer pair

231  ITS5a forward and ITS4 reverse. The presence of large amount of secondary metabolites,

232  polysaccharides and polyphenolic compounds in the plants of sub-family Caesalpinioidae,

233  hindered the isolation of pure nucleic acids. Therefore few samples had to be excluded from the

234  study after 3-4 initial amplification attempts that failed due to the presence of inhibitory

235  components. The present study generated 15 new sequences belonging to 15 different species of

236  *Cassia,   Senna,*   and   *Chamaecrista.*   The   sequences   were   submitted   to   NCBI

237  (www.ncbi.nlm.nih.gov/genbank/)  and  corresponding  GenBank  accession  numbers  were

238  obtained for each species. A total of 81 sequences corresponding to 18 different species of

239  *Cassia, Senna,* and *Chamaecrista* for ITS regions (ITS1 and ITS2) were obtained from NCBI

240  and included in the study (Table 1).   The ITS1 region had an aligned length of 315 bp

241  (Alignment S1) which was greater than that of ITS2 with 258 bp (Table 2; Alignment S2). The

242  combined region ITS1+2 showed an align length of 573 bp (Alignment S3) with 80.1 % of

243  pairwise identity (Table 2). The aligned ITS1 matrix consisted of 315 bp with 206 parsimony

244  sites. The number of variable sites was 210. The maximum intra-specific divergence was

245  observed among the individuals of *Senna siamea* with 0.023 PWG-distance while minimum

246  inter-specific distances were recorded between *Senna hirsuta* and *Senna occidentalis* with 0.039

247  PWG-distance. The species of genus *Chamaecrista* showed lowest K2P distances (Table 3).

248  Overall the summary statistics for DNA alignments and DNA sequences for the ITS dataset

249  evaluated in this study are summarized in Table 2 and Table 3 respectively.

250  **Genetic divergence and Barcoding gap**

251    The presence of DNA barcoding gap based on the concept of an  inter-specific distance being

252    larger than the intra-specific distance for a species, directly reveals the species discrimination

253    ability of candidate barcodes. In this study, the relative distribution of frequencies of K2P

254    distances for three ITS datasets using TaxonDNA software showed a significant pattern with the

255    inter-specific distance being higher and did not fully overlap with the intra-specific distance

256    resulting in the presence of an identified barcoding gap in the genera. The observed pattern of

257    ITS1, ITS2 and ITS1+2 results are presented in Figure 1. The mean intra- and inter-specific

258    genetic divergence based on PWG distances through MEGA, for ITS1 varied in the range from

259    0.023 to 0.000 and 0.033 to 1.185 respectively (Table 3).

260    **Species discrimination based on different analytical methods**

261    In accordance with the CBOL PWG-distance method, a favourable barcode should possess a

262    high inter-specific divergence to distinguish different species. The result obtained through the

263    different datasets showed significant pattern of inter-specific divergence, whereby ITS1 was

264    concluded to be the best among the candidates. The mean pairwise inter-specific distances were

265    found to be higher in comparison to intra-specific distances in all the barcodes, resulting in the

266    presence of a clear barcode gap. The distance distribution range of all inter- and intra-specific

267    distances for all markers are shown in Figure 2.

268         Compared with the PWG- distance method, the BM and BCM functions of TaxonDNA

269    showed the better discrimination success. All the three datasets presented same success rate of

270    species identification when BM was selected in comparison to BCM. The highest and same rate

271    of discriminatory power (81.6%) was observed for ITS1 on both BM and BCM functions. The

272    other two datasets; ITS2 and ITS1+2 datasets recovered 75.0% and 77.4% BM respectively

273    (Table 4).

274        The tree building methods for the evaluation of barcode sequences were estimated based

275    on the correct assignment of individuals forming a monophyletic clade (Figure 3 and Figure 1

276    Suppl.). Among the different phylogenetic methods, BI recovered the highest value for species

277    monophyly in all the datasets. While in the combination of ITS1+2, all the three methods viz.

278    NJ, ML and BI provided near similar topology, concluding 77.41% of individuals identified

279    correctly (Figure 4). The resulting bootstrap value lends support to our findings. Comparing the

280    potentiality of the ITS datasets and the phylogenetic algorithms employed, the highest

281    discriminatory power was observed when ITS1 was used alone, which successfully maintained

282    the genera (*Cassia, Senna,* and *Chamaecrista*) monophyly with few exceptions (Figure 5). The

283    coalescent and speciation tree priors intrinsically correlated the rate of evolution and time in

284    inferring genetic differences between species. It is interesting to conclude that all the species

285    from genera *Senna* and *Cassia* framed in two different clusters viz. Cluster I and II according to

286    traditional morphology. The phylogenetic tree presented a slight divergence in the clustering of

287    *Chamaecrista absus* accession obtained from GenBank which might be due to the mis-

288    identification of samples. Referring to the species relationships within genera; to some extent,

289    the phylogenetic relationships obtained were in consistent with the result obtained from the

290    traditional morphological classification method.  The clustering pattern of three different genera

291    *Cassia, Senna,* and *Chamaecrista* within the subtribe Cassinae based on the nuclear ribosomal

292    region ITS1, proved to be successful in comparison to the infrageneric clustering of taxa. The

293    clustering of *Senna tora*, *Senna uniflora* and *Senna obtusifolia* accessions based on molecular

294    algorithm of ITS1 complies with the morphological similarity occurs among them, while in

295    ITS2, *Senna uniflora* showed little divergence (Figure 3). Also we were not able to find out the

296    clear pattern of lineage of respective species within the genus at a molecular level, as according

297   to traditional taxonomy. Worthy to note here, that the resulting pattern within the individuals of

298   same species and high reliability value obtained for their nodes concludes the existence of

299   genetic similarity among them. Framing of *Senna occidentalis* and *Senna hirsuta* into the

300   individual cluster through ITS1, were in consistent with the key classification (Figure 3).

301   Besides, all the tree species belonging to genus *Cassia*, undertaken in this study framed

302   an individual cluster (Cluster II) according to their diversity there by concluding the importance

303   of molecular characterization in corroboration with morphological methods in biosystematics

304   study. The analysis conducted in subtribe Cassiinae with the tree based, similarity based and

305   distance based methods showed that BI phylogenetic method and BM similarity methods

306   outperformed the PWG- distance method when using these barcode loci (Figure 4).

307   **DISCUSSION**

308   **Discrimination success**

309   Hitherto several different analytical methods were framed for the assessment of the species

310   discrimination ability, which includes tree-based (NJ, MP, Bayesian), distance-based (PWG-

311   distance, p-distance, K2P-distance) and sequence similarity-based methods (Blast and

312   TaxonDNA), etc., and all of them show different discrimination power on the same data set

313   (*Little & Stevenson, 2007; Austerlitz et al., 2009; China Plant BOL Group, 2011; Sandionigi et*

314   *al., 2012*). In this study, sequence analysis of ITS datasets using Bayesian inference (BI) tree-

315   based method gave the highest species resolution based on the topology with the highest product

316   of posterior clade probabilities across all nodes followed by BM and BCM model of TaxonDNA,

317   which too presented equally efficient results either in single or combination of barcodes.

318   Similarly, patterned results have been obtained in different DNA barcoding studies in various

319    plant groups (*Yan et al., 2014; Giudicelli et al., 2015; Xu et al., 2015; Yan et al., 2015*). The

320    clustering algorithm of Bayesian framework provides a flexible way to model rate variation and

321    obtain reliable estimates of speciation times, provided the assumptions of the models be adequate

322    (*Drummond et al., 2012*).

323         The PWG-distance method based on simple pairwise matching recommended by CBOL

324    Plant Working Group as a universal and robust method for the assessment of clear barcoding gap

325    indicated the significance of ITS1, thereby highest number of variable and informative sites (210

326    and 206, respectively) were obtained. Moreover, the rate of species discrimination is equally

327    efficient when ITS1 and ITS2 are concatenated. These results were expected, considering the

328    complexity of the genera and directly reflected on the performance of ITS1 and ITS2 as barcode

329    markers in *Cassia, Senna,* and *Chamaecrista*. The possible reason behind the results might be the

330    inter-specific sharing of identical sequences or failure of conspecific individuals to group

331    together. Besides, many other aspects have also been reported for unclear barcoding gap such as

332    imperfect taxonomy, inter-specific hybridization, paralogy and incomplete lineage sorting (*Yan

333    et al., 2015*). However, ITS region has proved to be a suitable marker in authentication of *Cassia*

334    species in the commercial herbal market (*Seethapathy et al., 2014*). The strong identification

335    ability of nuclear region ITS have been verified in many complex groups (*Baldwin et al., 1995;

336    Alves et al., 2014; Wang et al., 2014; Giudicelli et al., 2015*). Therefore, we suggest that ITS1

337    itself could be the first option for DNA barcoding in subtribe Cassiinae, though ITS2 should not

338    be discarded.

339         Moreover, the differences among the three methods compared here, have their possible

340    cause in the theories behind their algorithms and the matter of comprehensive sampling. Thus the

341    comparison of species resolution between studies without consideration of the methods should be

342    avoided for one or the other reasons discussed, as species resolution is an important criterion for

343    assessment of robust barcodes.

**Biological implications of ITS based signalling in Cassiinae**

345         The corroboration of morphological, ecological, geographical, reproductive biology and

346    DNA sequence information paved the successful path for constructing robust taxonomy for

347    diverged plant taxa (*DeSalle et al., 2005; Fazekas et al., 2009; Hollingsworth et al., 2011*). The

348    ITS region appears to evolve more rapidly than coding regions in interpreting phylogenetic

349    relationships at lower taxonomic levels (Inter-generic and Inter-specific). Species discrimination

350    for the genera *Cassia, Senna* and *Chamaecrista* sampled in this study was high with the strong

351    identification ability of nuclear region ITS. All the three genera maintained the monophyly of the

352    clade either alone or in combination of barcoded loci. The resulting bootstrap value lends support

353    to our findings. To some extent, the divergence of species within the genus did not outperformed

354    as designated according to key taxonomy. The possible reasons behind the findings could be the

355    complexity of the genus with large number of highly polymorphic species which has been found

356    to devise greater interspecific variation (*Mohanty et al., 2010*). Sometimes interspecific

357    hybridization and gene introgression had accounted for the limited barcoding event at genus

358    level. Moreover genera *Cassia* and *Senna* accounts for high morphological complexity based on

359    species polymorphism, which have been reported in few studies in the past. Successful PCR

360    amplifications, sequencing strategy and alignment matrix obtained from the present study

361    provided further evidence to support the separation of species and genera. The robust

362    phylogenetic signalling of ITS region seems obvious in Cassinae. Although an earlier study

363    (excluding ITS) did not report any single novel region to differentiate the existing *Cassia* species

364    (*Purushothaman et al., 2014*), our findings provide the potentiality of the ITS region with data

365 support. The delineation of genera based on ITS regions provided a basic framework to have an

366 authentication prospect of correct species at the industrial level.

367 **CONCLUSIONS**

368 Our results show that ITS1 and ITS2 present all the desired characteristics of a DNA barcode for

369 the Cassiinae group examined in the present study. The high rate of PCR amplification and

370 sequencing success coupled with a potentially high rate of correctly assigned species among the

371 genera *Cassia, Senna,* and *Chamaecrista* conclude the discriminating capability of the nuclear

372 region ITS. However, till date, there has been much controversy over the ideal barcode for

373 plants. The previously advocated plastids regions have been used successfully in many barcoding

374 studies (*Kress & Erickson, 2007; CBOL Plant Working Group, 2009*). In many cases, the

375 potentiality of species discrimination based on the combination of ITS and plastid loci or ITS2

376 alone has been demonstrated in different plant groups (*Pang et al., 2010; Yang et al., 2012; Han*

377 *et al., 2013; Zhang et al., 2014*). The choice of ITS1 over ITS2, have been suggested recently in

378 the studied taxonomic group (*Wang et al., 2014*). Through our study, we concluded that ITS1

379 region should be used as a starting point to assign correct identification in the highly complex

380 genera *Cassia, Senna* and *Chamaecrista.*

381 **ACKNOWLEDGEMENTS**

385 **DNA Sequence Deposition**

386    The sequence data from this study has been submitted to the GenBank (NCBI) under Accession

387    Numbers KT279729.1– KT308097.1.

388    **Supplemental Information**

389    Figure 1 Suppl.: Phylogenetic consensus tree obtained for *Cassia*, *Senna,* and *Chamaecrista*

390    species based on nrITS datasets constructed using maximum likelihood algorithm.

391    AlignmentS1: The aligned sequences matrix of ITS1.

392    AlignmentS2: The aligned sequences matrix of ITS2.

393    AlignmentS3: Concatenated aligned sequences matrix of ITS1+2.

394    **References**

395    **Álvarez I, Wendel JF. 2003.** Ribosomal ITS sequences and plant phylogenetic inference.

396    *Molecular Phylogenetics and Evolution* 29:417-434 DOI 10.1016/S1055-7903(03)00208-

397    2.

398    **Alves TLS, Chauveau O, Eggers L, Souza-Chies TTD. 2014.** Species discrimination in

399    *Sisyrinchium* (Iridaceae): assessment of DNA barcodes in a taxonomically challenging

400    genus. *Molecular Ecology Resources* 14:324-335 DOI 10.1111/1755-0998.12182.

401    **Austerlitz F, David O, Schaeffer B, Bleakly K, Olteanu M, Leblois R, Veuile M, Laredo C.**

402    **2009.** DNA barcode analysis: a comparison of phylogenetic and statistical classification

403    methods. *BMC Bioinformatics* 10:S10 DOI 10.1186/1471-2105-10-S14-S10.

404    **Baldwin BG, Sanderson MJ, Porter JM, Wojciechowski MF, Campbell CS, Donoghue MJ.**

405    **1995.** The ITS region of nuclear ribosomal DNA: a valuable source of evidence on

406    angiosperm phylogeny. *Annals of the Missouri Botanical Garden* 82:247-277 DOI

407    10.2307/2399880.

408 **Bir SS, Kumari S. 1982.** Karyotipic studies in *Cassia* Linn. from India. *Proceedings of*
409     *the National Academy of Sciences, India, Section B: Biological Sciences* B48:397-404.

410 **Casiraghi M, Labra M, Ferri E, Galimberti A, deMattia F. 2010.** DNA barcoding: theoretical
411     aspects and practical applications. In: Nimis PL, Lebbe RV, eds. *Tools for Identifying*
412     *Biodiversity: Progress and Problems*. EUT Publishers, 269-273.

413 **CBoL Plant Working Group. 2009.** A DNA barcode for land plants. *Proceedings of the*
414     *National Academy of Sciences of the United States of America* 106:12794-12797 DOI
415     10.1073/pnas.0905845106.

416 **Chase MW, Salamin N, Wilkinson M, Dunwell JM, Kesanakurti RP, Haidar N, Savolainen**
417     **V. 2005.** Land plants and DNA barcodes: Short-term and long-term goals. *Philosophical*
418     *transactions of the Royal Society of London. Series B, Biological Sciences* 360:1889-
419     1895 DOI10.1098/rstb.2005.1720.

420 **Chase MW, Cowan RS, Hollingsworth PM, van den Berg C, Madrinan S, Petersen G,**
421     **Seberg O, Jorgsensen T, Cameron KM, Carine M. 2007.** A proposal for a standardised
422     protocol to barcode all land plants. *Taxon* 56:295-299.

423 **Chen S, Yao H, Han J, Liu C, Song J, Shi L, Zhu Y, Ma X, Gao T, Pang X, Luo K, Li Y, Li**
424     **X, Jia X, Lin Y, Leon C. 2010.** Validation of the ITS2 region as a novel DNA barcode
425     for identifying medicinal plant species. *PLoS One* 5:e8613 DOI
426     org/10.1371/journal.pone.0008613.

427 **China Plant BOL Group. 2011.** Comparative analysis of a large dataset indicates that internal
428     transcribed spacer (ITS) should be incorporates into the core barcode for seed plants.
429     *Proceedings of the National Academy of Sciences of the United States of America*
430     108:19641-19646 DOI 10.1073/pnas.1104551108.

431    **DeSalle R, Egan MG, Siddall M. 2005.** The unholy trinity: taxonomy, species delimitation and

432        DNA barcoding. *Philosophical transactions of the Royal Society of London. Series B,*

433        *Biological Sciences* 360:1905-1916 DOI 10.1098/RSTB.2005.1722.

434    **Drummond AJ, Suchard MA, Rambaut XD. 1973.** A Bayesian phylogenetics with BEAUti

435        and the BEAST 1.7. *Molecular Biology and Evolution* 29:1969-1973 DOI

436        10.1093/molbev/mss075.

437    **Fazekas AJ, Burgess KS, Kesanakurti PR, Graham SW, Newmaster SG, Husband BC,**

438        **Percy DM, Hajibabaei M, Barret SC. 2008.** Multiple multilocus DNA barcodes from

439        the plastid genome discriminate plant species equally well. *PLoS One* 3:e2802 DOI

440        org/10.1371/journal.pone.0002802.

441    **Fazekas AJ, Kesanakurti PR, Burgess KS, Perc DM, Graham SW, Barrett SC, Newmaster**

442        **SG, Hajibabaei M, Husband BC. 2009.** Are plant inherently harder to discriminate

443        than animal species using DNA barcoding markers? *Molecular Ecology Resources* 9:130-

444        139 DOI 10.1111/j.1755-0998.2009.02652.x.

445    **Felsenstein J. 1988.** Phylogenies from molecular sequences: inference and reliability. *Annual*

446        *Review of Genetics* 22:521-565 DOI 10.1146/annurev.ge.22.120188.002513.

447    **Gao T, Yao H, Song J, Liu C, Zhu Y, Ma X, Pang X, Xu H, Chen S. 2010.** Identification of

448        medicinal plants in the family Fabaceae using a potential DNA barcode ITS2. *Journal*

449        *of Ethnopharmacology* 130:116-121 DOI 10.1016/j.jep.2010.04.026.

450    **Giudicelli GC, Mäder G, Freitas de LB. 2015.** Efficiency of ITS Sequences for DNA

451        Barcoding in *Passiflora* (Passifloraceae). *International Journal of Molecular Sciences*

452        16:7289-7303 DOI 10.3390/ijms16047289.

453    **Gregory TR. 2005.** DNA barcoding does not compete with taxonomy. *Nature* 434:1067-1080

454        DOI 10.1038/4341067b.

455    **Hall TA. 1999.** BioEdit: a user-friendly biological sequence alignment editor and analysis

456        program for Windows 95/98/NT. *Nucleic Acids Symposium Series* 41:95-98.

457    **Han J, Zhu Y, Chen X Liao B, Yao H, Song J, Chen S, Meng F. 2013.** The short ITS2

458        sequence serves as an efficient taxonomic sequence tag in comparison with the full-

459        length ITS. *BioMed Research International* 2013:741-476 DOI g/10.1155/2013/741476.

460    **Hebert PDN, Cywinska A, Ball SL, deWaard JR. 2003.** Biological identification through

461        DNA barcodes. *Proceedings of the Royal Society B: Biological Sciences* 270:313-321

462        DOI 10.1098/rspb.2002.2218.

463    **Hollingsworth ML, Andra Clark A, Forrest LL, Richardson J, Pennington RT, Long DG,**

464        **Cowan R, Chase MW, Gaudeul M, Hollingsworth PM. 2009.** Selecting barcoding loci

465        for plants: Evaluation of seven candidate loci with species-level sampling in three

466        divergent groups of land plants. *Molecular Ecology Resources* 9:439-457 DOI

467        10.1111/j.1755-0998.2008.02439.

468    **Hollingsworth PM, Graham SW, Little DP. 2011.** Choosing and using a plant DNA barcode.

469        *Plos One* 6:e19254 DOI org/10.1371/journal.pone.0019254.

470    **Huson DH, Scornavacca C. 2012.** Dendroscope 3: An interactive tool for rooted phylogenetic

471        trees and networks. *Systematic Biology* 61:1061-1067 DOI 10.1093/sysbio/sys062.

472    **Irwin HS, Barneby RC. 1981.** Tribe Cassieae Bronn. In: Polhill RM, Raven PH, eds. Recent

473        advances in legume systematics. Kew: Royal Botanic Garden. 97-106.

474    **Kress WJ, Wurdack KJ, Zimmer EA, Weigt LA, Janzen DH. 2005.** Use of DNA barcodes to

475        identify flowering plants. *Proceedings of the National Academy of Sciences of the United*

476        *States of America* 102:8369-8374 DOI 10.1073/pnas.0503123102.

477    **Kress WJ, Erickson DL. 2007.** A two locus global DNA barcode for land plants: The coding

478        *rbcL* gene complements the noncoding *trnH-psbA* spacer region. *Plos One* 2:e508 DOI

479        org/10.1371/journal.pone.0000508.

480    **Khanuja SPS, Shasany AK, Darokar MP, Kumar S. 1999.** Rapid isolation of DNA from dry

481        and fresh samples of plants producing large amounts of secondary metabolites and

482        essential oils. *Plant Molecular Biology Reporter* 17:1-7 DOI 10.1023/A:

483        1007528101452.

484    **Kimura M. 1980.** A simple method for estimating evolutionary rates of base substitutions

485        through comparative studies of nucleotide sequences. *Journal of Molecular Evolution*

486        16:111-120.

487    **Lahaye R, Van der Bank M, Bogarin D, Warner J, Pupulin F, Gigot G, Maurin O, Duthoit**

488        **S, Barraclough TG, Savolainen V. 2008.** DNA barcoding the floras of biodiversity

489        hotspots. *Proceedings of the National Academy of Sciences of the United States of*

490        *America* 105:2923-2928 DOI 10.1073/pnas.0709936105.

491    **Little DP, Stevenson DW. 2007.** A comparison of algorithms for the identification of specimens

492        using DNA barcodes: examples from gymnosperms. *Cladistics* 23:1-21 DOI

493        10.1111/j.1096-0031.2006.00126.

494    **Li X, Yang Y, Henry RJ, Rosseto M, Wang Y, Chen S. 2015.** Plant DNA barcoding: from

495        gene to genome. *Biological Reviews* 90:157-166 DOI 10.1111/brv.12104.

496    **Meier R, Shiyang K, Vaidya G, Ng PK. 2006.** DNA barcoding and taxonomy in Diptera: a tale

497    of high intraspecific variability and low identification success. *Systematic Biology*

498    55:715-728 DOI 10.1080/10635150600969864.

499    **Meier R, Zhang G, Ali F. 2008.** The use of mean instead of smallest interspecific distances

500    exaggerates the size of the "Barcoding Gap" and leads to misidentification. *Systematic*

501    *Biology* 57:809-813 DOI 10.1080/10635150802406343.

502    **Meyer CP, Paulay G. 2005.** DNA barcoding: error rates based on comprehensive sampling.

503    *PLoS Biology* 3:e422 DOI org/10.1371/journal.pbio.0030422.

504    **Mishra P, Kumar A, Nagireddy A, Mani D, Shukla AK, Tiwari R, Sundaresan V. 2016.**

505    DNA barcoding: an efficient tool to overcome authentication challenges in the herbal

506    market. *Plant Biotechnology Journal* 14:8-21 DOI 10.1111/pbi.12419.

507    **Mohanty S, Das AB, Gosh N, Panda BB, Smithe DW. 2010.** Genetic diversity of 28 wild

508    species of fodder legume *Cassia* using RAPD, ISSR and SSR markers: a novel breeding

509    strategy. *Journal of Biological Research* 2:44-55 DOI

510    **Nilsson RH, Ryberg M, Kristiansson E, Abarenkov K, Larsson KH, Koljalg U. 2006.**

511    Taxonomic Reliability of DNA Sequences in Public Sequence Databases: A Fungal

512    Perspective. *PLoS One* 1:e59 DOI org/10.1371/journal.pone.0000059.

513    **Pang X, Song J, Zhu Y, Xie C, Chen S. 2010.** Using DNA barcoding to identify species within

514    Euphorbiaceae. *Planta Medica* 76:1784-1786 DOI 10.1055/s-0030-1249806.

515    **Purushothaman N, Newmaster SG, Ragupathy S, Stalin S, Suresh D, Arunraj DR,**

516    **Gnanasekaran G, Vassou SL, Narasimhan D, Parani M. 2014.** A tiered barcode

517    authentication tool to differentiate medicinal *Cassia* species in India.

518    *Genetics and Molecular Research* 13:2959-2968 DOI 10.4238/2014.April.16.4.

519    **Rambaut A, Suchard MA, Xie D, Drummond AJ. 2014.** Tracer v1.6, Available

520         from http://beast.bio.ed.ac.uk/Tracer.

521    **Sandionigi A, Galimberti A, Labra M, Ferri E, Panunzi E, deMattia F, Casiraghi M. 2012.**

522         Analytical approaches for DNA barcoding data-how to find a way for plants? *Plant*

523         *Biosystems* 146:805-813 DOI 10.1080/11263504.2012.740084.

524    **Savolainen V, Cowan RS, Vogler AP, Roderick GK, Lane R. 2005.** Towards writing the

525         encyclopaedia of life: An introduction to DNA barcoding. *Philosophical transactions of*

526         *the Royal Society of London. Series B, Biological Sciences* 360:1805-1811 DOI

527         10.1098/rstb.2005.1730.

528    **Seethapathy GS, Ganesh D, Santhosh Kumar JU, Senthilkumar U, Newmaster SG,**

529         **Ragupathy S, Shaanker RU, Ravikanth G. 2014.** Assessing product adulteration in

530         natural health products for laxative yielding plants, *Cassia*, *Senna*, and *Chamaecrista* in

531         Southern India using DNA barcoding. *International Journal of Legal Medicine* DOI

532         10.1007/s00414-014-1120-z.

533    **Singh V. 2001.** *Monograph on Indian subtribe Cassinae (Cesalpiniaceae)*. Scientific Publisher:

534         India.

535    **Swofford DL. 2003.** PAUP*: Phylogenetic analysis using parsimony (* and other methods),

536         version 4.0b10. Sunderland: Sinauer.

537    **Tamura K, Stecher G, Peterson D, Filipski A, Kumar S. 2013.** MEGA6: Molecular

538         evolutionary genetics analysis version 6.0. *Molecular Biology and Evolution* 30:2725-

539         2729 DOI 10.1093/molbev/mst197/

540    **Vijayan K, Tsou CH. 2010.** DNA barcoding in plants: taxonomy in a new perspective. *Current*

541         *Science India*. 99:1530-1541 DOI

542  **Wang XC, Liu C, Huang L, Bengtsson-Palme J, Chen H, Zhang JH, Cai D, Li JQ. 2014.**

543      ITS1: A DNA barcode better than ITS2 in eukaryotes? *Molecular Ecology Resources*

544      DOI 10.1111/1755-0998.12325.

545  **Xu S, Li D, Li J, Xiang X, Jin W, Huang W, Xiaohua J, Huang L. 2015.** Evaluation of the

546      DNA Barcodes in *Dendrobium* (Orchidaceae) from Mainland Asia. *PLoS One*

547      10:e0115168 DOI rg/10.1371/journal.pone.0115168.

548  **Yang JB, Wang YP, Möller M, Gao LM, Wu D. 2012.** Applying plant DNA barcodes to

549      identify species of *Parnassia* (Parnacciaceae). *Molecular Ecology Resources* 12:267-275

550      DOI 10.1111/j.1755-0998.2011.03095.

551  **Yan LJ, Liu J, Moller M, Zhang L, Zhang XM, Li DZ, Gao LM. 2014.** DNA barcoding of

552      *Rhododendron* (Ericaceae), the largest Chinese plant genus in biodiversity hotspots of the

553      Himalaya-Hengduan Mountains. *Molecular Ecology Resources* DOI 10.1111/1755-

554      0998.12353.

555  **Yan HF, Liu YJ, Xie XF, Zhang CY, Hu CM, Hao G, Ge XJ. 2015.** DNA barcoding

556      evaluation and its taxonomic implications in the species rich genus *Primula* L. in China.

557      *PLoS One* 10:e0122903 DOI org/10.1371/journal.pone.0122903.

558  **Yadav P, Koul KK, Srivastava N, Mendki MJ, Bhagyawant SS. 2016.** ITS-PCR sequencing

559      approach deciphers molecular phylogeny in chickpea, *Plant Biosystems - An*

560      *International Journal Dealing with all Aspects of Plant Biology* DOI

561      10.1080/11263504.2016.1179694.

562  **Zhang D, Duan L, Zhou N. 2014.** Application of DNA barcoding in *Roscoea* (Zingiberaceae)

563      and a primary discussion on taxonomic status of *Roscoea cautleoides* var. *Pubescens*.

564      *Biochem Systematics and Ecology* 52:14-19 DOI 10.1016/j.bse.2013.10.004.

566    **Figure legends**

567    **Figure 1 Relative abundance of intra- and inter-specific Kimura-2-Parameter pairwise**

568    **distance based on TaxonDNA methods considering nrITS dataset in genera *Cassia, Senna,***

569    **and *Chamaecrista*.**

570    **Figure 2 Relative distributions of intra- and inter-specific distances based on PWG-**

571    **distance based methods for the three nrITS datasets in Cassiinae.** x axes relate to Kimura 2-

572    parameter (K2P) distances arranged in intervals, and the y axes correspond to the frequency

573    distribution.

574    **Figure 3 Phylogenetic consensus tree obtained for *Cassia*, *Senna,* and *Chamaecrista* species**

575    **based on nrITS datasets constructed using bayesian inference algorithm.** Representatives

576    from individual species are abbreviated based on corresponding taxon.

577    **Figure 4 Species discrimination rates of nrITS datasets based on different methods in**

578    **Cassiinae.** ITS1 barcode in conjunction with the bayesian inference analysis of hierarchical tree-

579    based method met the objectives of DNA barcoding.

580    **Figure 5 Evolutionary relationships in genera *Cassia*, *Senna,* and *Chamaecrista* based on**

581    **nrITS barcode constructed using bayesian inference algorithm.** Taxon names are abbreviated

582    (see Table 1).

# Table 1(on next page)

Passport sheet for the samples undertaken.

Sample details with GenBank accession numbers of all the samples of *Cassia, Senna, and Chamaecrista* used in this study. Accessions numbers marked in bold represent lab-generated sequences from the present study.

1  **Table 1 Sample details with GenBank accession numbers of all the samples of *Cassia,***
2  ***Senna, and Chamaecrista* used in this study. Accessions numbers marked in bold represent**
3  **lab-generated sequences from the present study.**

| Taxon | Region | Collection Site | Voucher Number (No.) | GenBank (NCBI) Accessions No. |
|---|---|---|---|---|
| *Chamaecrista absus* | ITS | Tirunelveli, Tamil Nadu | CIMAP-C010 | **KT279729.1** |
| *Chamaecrista absus* | ITS2 | GenBank | GenBank | FJ009832.1 |
| *Chamaecrista absus* | ITS | GenBank | GenBank | KC817015.1 |
| *Chamaecrista absus* | ITS2 | GenBank | GenBank | FJ009832.1 |
| *Chamaecrista nigricans* | ITS | Tuticorin, Tamil Nadu | CIMAP-C011 | **KT279731.1** |
| *Chamaecrista nigricans* | ITS2 | GenBank | GenBank | JQ301845.1 |
| *Chamaecrista nigricans* | ITS | Tuticorin, Tamil Nadu | CIMAP-C011 | **KT279731.1** |
| *Chamaecrista nigricans* | ITS2 | GenBank | GenBank | JQ301845.1 |
| *Senna uniflora* | ITS | Tirunelveli, Tamil Nadu | CIMAP-C012 | **KT279730.1** |
| *Senna uniflora* | ITS | GenBank | GenBank | KJ605909.1 |
| *Senna uniflora* | ITS | GenBank | GenBank | KJ605897.1 |
| *Senna italica* | ITS | Tuticorin, Tamil Nadu | CIMAP-C013 | **KT279732.1** |
| *Senna italica* | ITS | GenBank | GenBank | KJ004293.1 |
| *Senna italica* | ITS | GenBank | GenBank | KF815503.1 |
| *Senna hirsuta* | ITS | Tirunelveli, Tamil Nadu | CIMAP-C014 | **KT279733.1** |
| *Senna hirsuta* | ITS | GenBank | GenBank | KJ605904.1 |
| *Cassia fistula* | ITS2 | GenBank | GenBank | JQ301830.1 |
| *Senna hirsuta* | ITS | GenBank | GenBank | KJ605905.1 |
| *Senna hirsuta* | ITS2 | GenBank | GenBank | KJ605904.1 |
| *Senna alata* | ITS | Kukrail, Lucknow | CIMAP-C015 | **KT308089.1** |
| *Senna alata* | ITS | GenBank | GenBank | KJ638414.1 |
| *Senna alata* | ITS | GenBank | GenBank | KJ638413.1 |
| *Senna sulfurea* | ITS | Raebareli, Lucknow | CIMAP-C016 | **KT308090.1** |
| *Senna sulfurea* | ITS2 | GenBank | GenBank | JQ301833.1 |
| *Senna siamea* | ITS | CIMAP, Bangalore | CIMAP-C017 | **KT308091.1** |
| *Senna siamea* | ITS | GenBank | GenBank | KC984644.1 |
| *Senna siamea* | ITS | GenBank | GenBank | KJ638421.1 |
| *Senna siamea* | ITS2 | GenBank | GenBank | JQ301842.1 |
| *Senna obtusifolia* | ITS | Raebareli, Lucknow | CIMAP-C018 | **KT308092.1** |
| *Senna obtusifolia* | ITS | GenBank | GenBank | GU175319.1 |
| *Senna occidentalis* | ITS | Frlht, Bangalore | CIMAP-C019 | **KT308093.1** |
| *Senna occidentalis* | ITS | GenBank | GenBank | KJ638419.1 |
| *Senna occidentalis* | ITS | GenBank | GenBank | KP092706.1 |
| *Senna occidentalis* | ITS2 | GenBank | GenBank | KJ638419.1 |
| *Senna occidentalis* | ITS2 | GenBank | GenBank | KP092706.1 |
| *Senna pallida* | ITS | Raebareli, Lucknow | CIMAP-C020 | **KT308095.1** |
| *Cassia fistula* | ITS2 | GenBank | GenBank | JQ301830.1 |
| *Senna pallida* | ITS2 | GenBank | GenBank | JQ301829.1 |
| *Senna auriculata* | ITS | Frlht, Bangalore | CIMAP-C021 | **KT308096.1** |
| *Senna auriculata* | ITS | GenBank | GenBank | KJ638417.1 |
| *Senna auriculata* | ITS2 | GenBank | GenBank | JQ301838.1 |

| | | | | |
|---|---|---|---|---|
| *Senna auriculata* | ITS | GenBank | GenBank | KJ638416.1 |
| *Senna alexandrina* | ITS | CIMAP, Lucknow | CIMAP-C022 | **KT308097.1** |
| *Senna alexandrina* | ITS | GenBank | GenBank | KF815491.1 |
| *Senna alexandrina* | ITS2 | GenBank | GenBank | JQ301846.1 |
| *Senna alexandrina* | ITS2 | GenBank | GenBank | JQ301846.1 |
| *Senna surattensis* | ITS | GenBank | GenBank | KJ638427.1 |
| *Senna surattensis* | ITS | GenBank | GenBank | KJ605903.1 |
| *Senna surattensis* | ITS | GenBank | GenBank | KJ605902.1 |
| *Senna surattensis* | ITS2 | GenBank | GenBank | KJ638427.1 |
| *Senna tora* | ITS | GenBank | GenBank | KJ638426.1 |
| *Senna siamea* | ITS2 | GenBank | GenBank | JQ301842.1 |
| *Senna tora* | ITS | GenBank | GenBank | KJ638425.1 |
| *Senna tora* | ITS | GenBank | GenBank | KJ638424.1 |
| *Senna tora* | ITS2 | GenBank | GenBank | KJ638426.1 |
| *Senna tora* | ITS2 | GenBank | GenBank | KJ638425.1 |
| *Senna tora* | ITS2 | GenBank | GenBank | KJ638424.1 |
| *Cassia roxburghii* | ITS | GenBank | GenBank | JX856435.1 |
| *Cassia roxburghii* | ITS2 | GenBank | GenBank | JQ301841.1 |
| *Cassia javanica* | ITS | Raebareli, Lucknow | CIMAP-C023 | **KT338798.1** |
| *Cassia javanica* | ITS | GenBank | GenBank | FJ009821.1 |
| *Cassia javanica* | ITS2 | GenBank | GenBank | JQ301831.1 |
| *Cassia javanica* | ITS | GenBank | GenBank | FJ980413.1 |
| *Cassia javanica* | ITS2 | GenBank | GenBank | JQ301831.1 |
| *Cassia fistula* | ITS | SCAD, Tirunelveli | CIMAP-C024 | **KT308094.1** |
| *Cassia fistula* | ITS | GenBank | GenBank | JX856431.1 |
| *Cassia fistula* | ITS | GenBank | GenBank | JX856430.1 |
| *Cassia fistula* | ITS2 | GenBank | GenBank | JQ301830.1 |
| *Senna surattensis* | ITS2 | GenBank | GenBank | KJ638427.1 |
| *Senna surattensis* | ITS2 | GenBank | GenBank | KJ638427.1 |
| *Senna pallida* | ITS | Raebareli, Lucknow | CIMAP-C020 | **KT308095.1** |
| *Senna pallida* | ITS2 | GenBank | GenBank | JQ301829.1 |
| *Senna auriculata* | ITS2 | GenBank | GenBank | JQ301838.1 |
| *Senna auriculata* | ITS2 | GenBank | GenBank | JQ301838.1 |
| *Senna hirsuta* | ITS2 | GenBank | GenBank | KJ605904.1 |
| *Senna hirsuta* | ITS2 | GenBank | GenBank | KJ605904.1 |
| *Senna siamea* | ITS2 | GenBank | GenBank | JQ301842.1 |
| *Cassia javanica* | ITS2 | GenBank | GenBank | JQ301831.1 |
| *Cassia javanica* | ITS | GenBank | GenBank | FJ009821.1 |
| *Cassia roxburghii* | ITS | GenBank | GenBank | JX856435.1 |
| *Cassia roxburghii* | ITS2 | GenBank | GenBank | JQ301841.1 |

4

**Table 2**(on next page)

Summary for DNA alignments.

Summary statistics for DNA alignments.

1    **Table 2 Summary statistics for DNA alignments.**

| Alignments | Region | Residual length | G+C (%) | Identical sites (%) | Pairwise identity (%) |
|---|---|---|---|---|---|
| Alignment S1 | ITS1 | 315 | 57.0 % | 26.3 % | 82.15 % |
| Alignment S2 | ITS2 | 258 | 63.9 % | 35.8 % | 77.20 % |
| Alignment S1+2 | ITS1+2 | 573 | 60.1 % | 30.8 % | 80.10 % |

2    **Notes.**

3    *Residual length,* the length of the complete alignment, counting portions excluded from analysis; *G+C*, the G + C
4    content of the complete (total length) alignment; *Identical sites*, the % of columns in the alignment for which all
5    sequences are identical; *Pairwise identity*, the % of pairwise residues that are identical in the alignments, including
6    gap versus non-gap residues, but excluding gap vs. gap residues.

7

**Table 3**(on next page)

Summary of sequence characteristics

Summary of sequence characteristics of the barcode candidates and their combinations
analysed in this study.

1 **Table 3 Summary of sequence characteristics of the barcode candidates and their**
2 **combinations analysed in this study.**

| Characters | ITS1 | ITS2 | ITS1+2 |
|---|---|---|---|
| Aligned length (bp) | 315 | 258 | 573 |
| Average intra-distance | 0.01% | 0.03% | 0.01% |
| Average inter-distance | 0.24% | 0.25% | 0.17% |
| Average theta (e) | 0.27% | 0.26% | 0.18% |
| Coalescent depth | 0.02% | 0.38% | 0.17% |
| Proportion of variable sites | 66.66% | 60.24% | 46.53% |
| Proportion of parsimony sites | 65.39% | 47.54% | 43.64% |

3

**Table 4**(on next page)

Identification success rates based on analysis function of TaxonDNA software

Identification success rates based on analysis of the 'Best match', 'Best close match' and 'All species barcodes' function of TaxonDNA software for each ITS dataset.

1 **Table 4 Identification success rates based on analysis of the 'Best match', 'Best close**
2 **match' and 'All species barcodes' function of TaxonDNA software for each ITS dataset.**

| Region | Best match | | | Best close match | | | All species barcodes | | |
|---|---|---|---|---|---|---|---|---|---|
| | Correct (%) | Ambiguous (%) | Incorrect (%) | Correct (%) | Ambiguous (%) | Incorrect (%) | Correct (%) | Ambiguous (%) | Incorrect (%) |
| ITS1 | 81.63 | 8.16 | 10.2 | 81.63 | 8.16 | 10.2 | 30.61 | 63.26 | 6.12 |
| ITS2 | 75.0 | 0 | 25.0 | 75.0 | 0 | 25.0 | 33.33 | 62.5 | 4.16 |
| ITS1+2 | 77.41 | 19.35 | 3.22 | 77.41 | 19.35 | 3.22 | 19.35 | 77.41 | 3.22 |

3

# Figure 1

. Pairwise distance based on K2P method.

Relative abundance of intra- and inter-specific Kimura-2-Parameter pairwise distance based on TaxonDNA methods considering nrITS dataset in genera *Cassia*, *Senna*, and *Chamaecrista.*
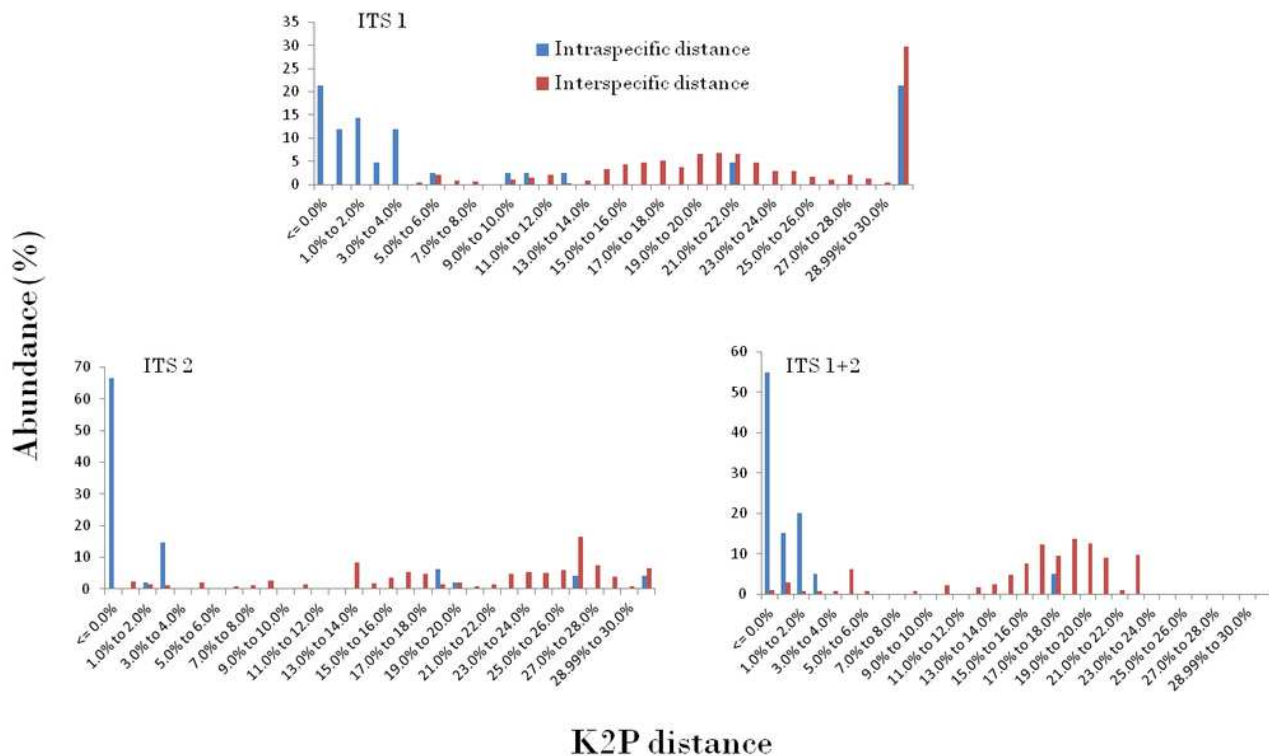
# Figure 2

Evaluation of DNA barcoding Gap

**Relative distributions of intra- and inter-specific distances based on PWG-distance based methods for the three nrITS datasets in Cassiinae.** x axes relate to Kimura 2-parameter (K2P) distances arranged in intervals, and the y axes correspond to the frequency distribution.
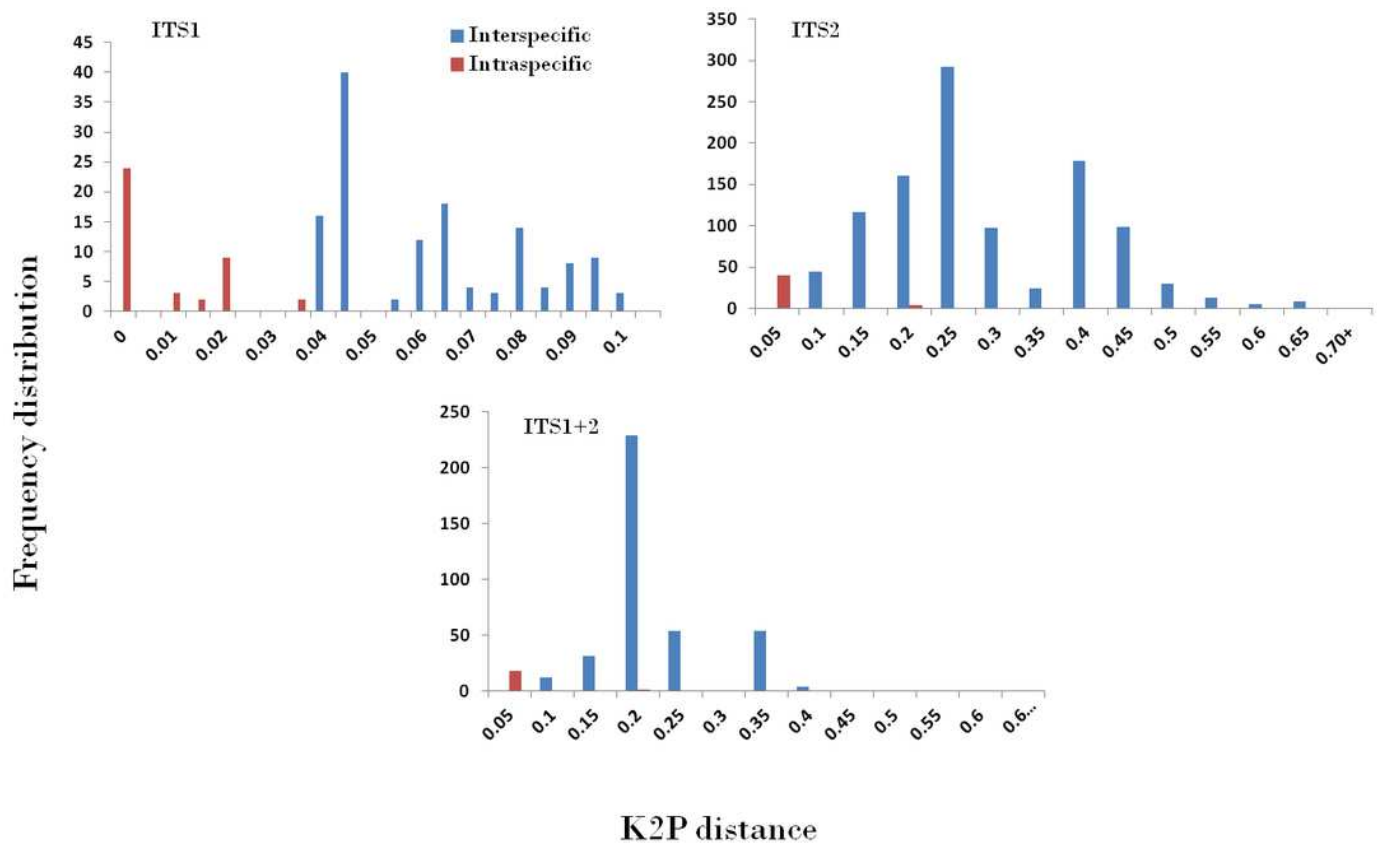
# Figure 3

Phylogenetic consensus tree constructed using bayesian inference algorithm.

**Phylogenetic consensus tree obtained for *Cassia*, *Senna*, and *Chamaecrista* species based on nrITS datasets constructed using bayesian inference algorithm.**

Representatives from individual species are abbreviated based on corresponding taxon.
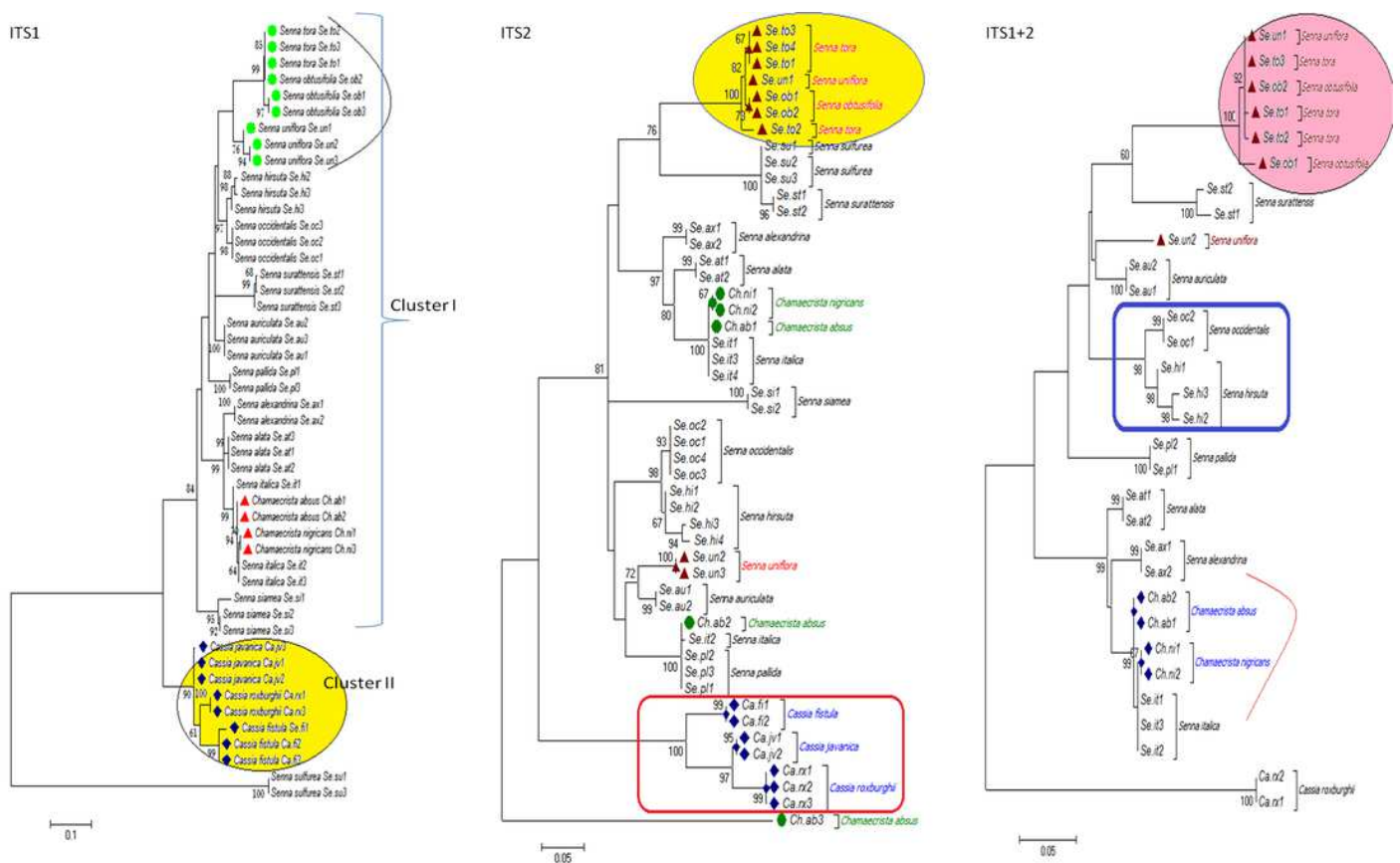
# Figure 4

Comparison of species discrimination rates

**Species discrimination rates of nrITS datasets based on different methods in Cassiinae.** ITS1 barcode in conjunction with the bayesian inference analysis of hierarchical tree-based method met the objectives of DNA barcoding.
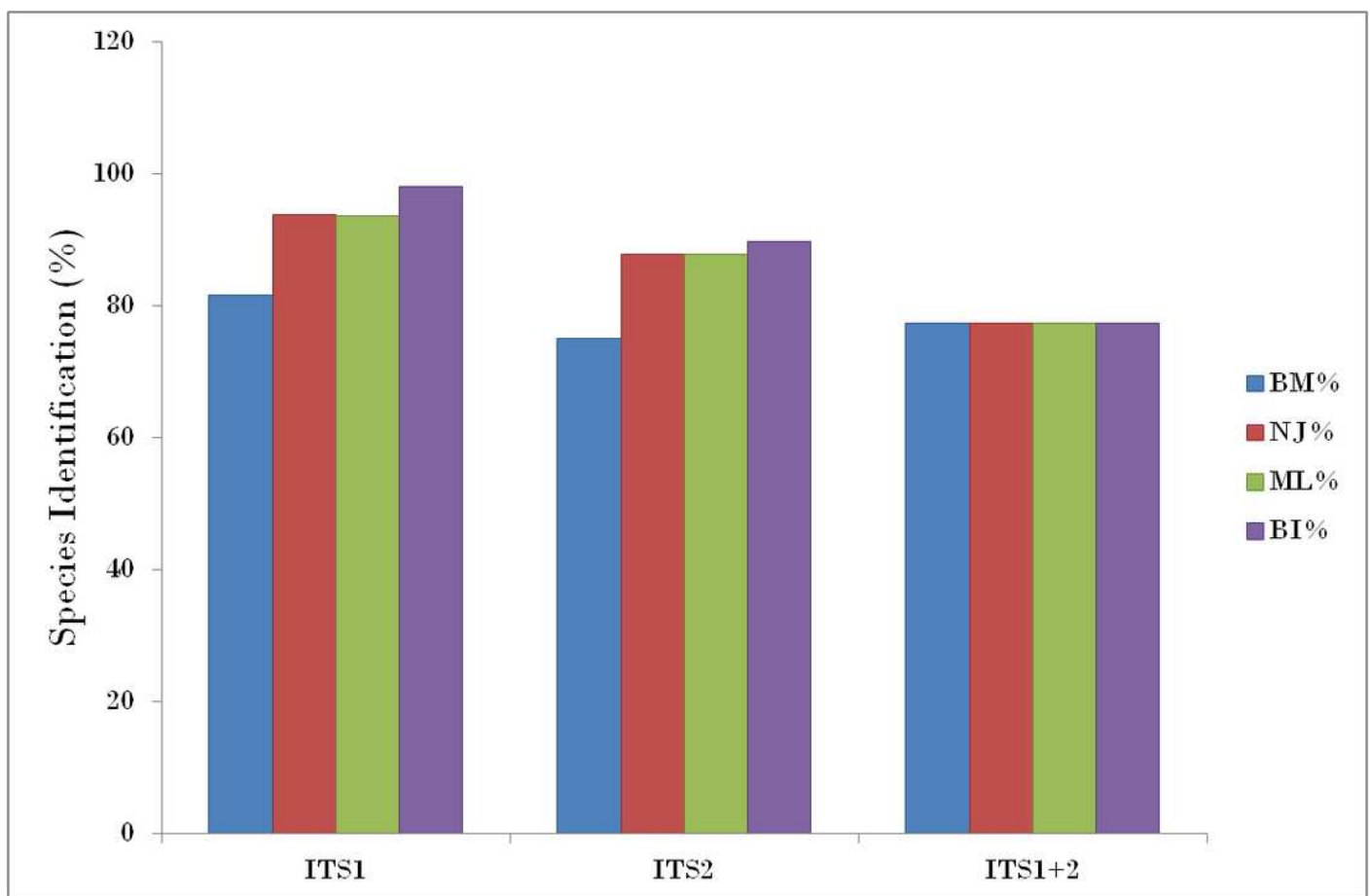
# Figure 5

Evolutionary relationships in Cassinae.

**Evolutionary relationships in genera *Cassia*, *Senna,* and *Chamaecrista* based on nrITS barcode constructed using bayesian inference algorithm.** Taxon names are abbreviated (see Table 1).