1    **Optimal markers for the identification of *Colletotrichum* species**

2

3    Willie Anderson dos Santos Vieira[a]

4    Priscila Alves Bezerra[a]

5    Anthony Carlos da Silva[a]

6    Josiene Silva Veloso[a]

7    Marcos Paz Saraiva Câmara[a]

8    Vinson Patrick Doyle[b]

9

10   [a] Universidade Federal Rural de Pernambuco – UFRPE. Departamento de

11   Agronomia. Rua Manuel de Medeiros, s/n - Dois Irmãos, Recife – Pernambuco,

12   52171-900.

13   [b] Department of Plant Pathology and Crop Physiology, Louisiana State University –

14   LSU, AgCenter, Baton Rouge, Louisiana, United States of America, 70808.

15

16   E-mail (in the author's order): andersonvieira12@gmail.com;

17   priscilaalvesbezerra@gmail.com; anthonycarlos17@hotmail.com;

18   josieneveloso@yahoo.com.br; marcos.camara@ufrpe.br; vdoyle@agcenter.lsu.edu.

19

20   Corresponding author: Willie Anderson dos Santos Vieira

21

22   **ABSTRACT**

23

24      *Colletotrichum* is among the most important genera of fungal plant pathogens.

25   Molecular phylogenetic studies over the last decade have resulted in a much better

26   understanding of the evolutionary relationships and species boundaries within the

27   genus. There are now approximately 200 species accepted, most of which are

28   distributed among 13 species complexes. Given their prominence on agricultural

29   crops around the world, rapid identification of a large collection of *Colletotrichum*

30   isolates is routinely needed by plant pathologists, regulatory officials, and fungal

31   biologists. However, there is no agreement on the best molecular markers to

32   discriminate species in each species complex. Here we calculate the barcode gap

33   distance and intra/inter-specific distance overlap to evaluate each of the most

34   commonly applied molecular markers for their utility as a barcode for species

35   identification. Glyceraldehyde-3-phosphate dehydrogenase (GAPDH), histone-3

36   (HIS3), DNA lyase (APN2),  intergenic spacer between DNA lyase and the mating-

37   type locus *MAT*1-2-1 (APN2/MAT-IGS), and intergenic spacer between GAPDH and

38   a hypothetical protein (GAP2-IGS) have the properties of good barcodes, whereas

39   sequences of actin (ACT), chitin synthase (CHS-1) and nuclear rDNA internal

40   transcribed spacers (nrITS) are not able to distinguish most species. Finally, we

41   assessed the utility of these markers for phylogenetic studies using phylogenetic

42   informativeness profiling, the genealogical sorting index (GSI), and Bayesian

43   concordance analyses (BCA). Although GAPDH, HIS3 and β-tubulin (TUB2) were

44   frequently among the best markers, there was not a single set of markers that were

45   best for all species complexes. Eliminating markers with low phylogenetic signal

46   tends to decrease uncertainty in the topology, regardless of species complex, and

47    leads to a larger proportion of markers that support each lineage in the Bayesian

48    concordance analyses. Finally, we reconstruct the phylogeny of each species

49    complex using a minimal set of phylogenetic markers with the strongest phylogenetic

50    signal and find the majority of species are strongly supported as monophyletic.

51

52    **KEYWORDS**

53    Accuracy

54    Anthracnose

55    Barcoding

56    Phylogenetic informativeness

57    Standardization

58

59

60    **1. INTRODUCTION**

61    *Colletotrichum* is among the largest groups of phytopathogenic fungi and

62    includes the causal agents of anthracnose and other diseases on seeds, stems,

63    leaves and fruits of important temperate and tropical crops (Cai et al., 2009, Cannon

64    et al., 2012). It is also among the most common genera of endophytic fungi, fungi

65    that live within plant organs without producing any symptoms of disease (Cannon et

66    al., 2012). Due to its economic and scientific importance, *Colletotrichum* was ranked

67    as the eighth most important phytopathogenic fungus in the world by plant

68    pathologists (Dean et al., 2012).

69    Species identification is necessary to understand disease epidemiology and

70    develop strategies to control the disease successfully (Cai et al., 2009). However,

71    *Colletotrichum* taxonomy and systematics has been a challenge since the genus was

72    introduced by Corda (1831). *Colletotrichum* species were historically circumscribed

73    on the basis of phenotypic features and a strong emphasis on the host species from

74    which the specimens were isolated under the assumption of host specificity, which

75    led to more than 900 species being recognized until revisionary work more than a

76    century after its introduction (von Arx, 1957; Sutton, 1980). *Colletotrichum*

77    identification based on morphological characters is problematic due to plasticity and

78    variation induced by experimental conditions (Vieira et al. 2017), and all life stages

79    are not frequently produced in culture (Samarakoon et al., 2018). The absence of

80    stable phenotypic characters has limited our understanding of phylogenetic

81    relationships within *Colletotrichum* and made the recognition of species boundaries

82    unreliable and confusing (Cai et al., 2009). To address this problem, Cai et al. (2009)

83    proposed a guideline for *Colletotrichum* species recognition based on a polyphasic

84    approach, which comprises the use of cultural, morphological, physiological and

85    pathogenicity characters in combination with phylogenetic analysis of nucleic acid

86    sequences.

87        The earliest phylogenetic studies of *Colletotrichum* using DNA sequences

88    were published by Mills et al. (1992) and Sreenivasaprasad et al. (1992).

89    Polymorphisms in the ITS1 region of the nrDNA were used to distinguish

90    *Colletotrichum* species. However, while the nrITS region is the most widely

91    sequenced region and has been chosen as the barcode locus for the Fungi, the utility

92    of this region is limited for systematic studies in *Colletotrichum*. Species diversity is

93    usually underestimated when based on nrITS sequences alone (Crouch et al.,

94    2009a) and it has been demonstrated to have little phylogenetic utility (Doyle et al.

95    2013; Vieira et al. 2017).  However, several additional markers have been applied for

96    multilocus phylogenetic inference to resolve the boundaries of cryptic species in the

97    genus (Cai et al. 2009; Damm et al. 2009; Doyle et al., 2013; Hyde et al., 2009; Lima

98    et al., 2013; Liu et al., 2016; Samarakoon et al., 2018; Veloso et al., 2018; Vieira et

99    al., 2014, 2017).

100        According to the most recent synopsis of the genus published in 2017, 188

101    *Colletotrichum* species have been described and incorporated into molecular

102    phylogenies. Among these species, 164 were distributed among 11 species

103    complexes and an additional 24 species had not been assigned to a species complex

104    (Marin-Felix et al., 2017). Additional species were recently described and three

105    additional clades were declared to represent new species complexes (Cao et al.,

106    2018; Damm et al., 2019; Samarakoon et al., 2018). Due to its global distribution and

107    ecological and economic importance, research groups around the world are working

108    concomitantly to address regional diversity. However, the set of phylogenetic

109    markers used to discriminate species is variable by species complex and no standard

110    set of markers has been adopted based on objective criteria (Marin-Felix et al.,

111    2017), making it difficult to combine data from disparate research groups and reliably

112    infer phylogenies and delimit species boundaries. Currently, thirteen different

113    molecular markers are commonly sequenced among the various *Colletotrichum*

114    species complexes: actin (ACT), DNA lyase (APN2), intergenic spacer between DNA

115    lyase and the mating-type locus *MAT*1-2-1 (APN2/MAT-IGS), calmodulin (CAL),

116    chitin synthase (CHS-1), glyceraldehyde-3-phosphate dehydrogenase (GAPDH),

117    intergenic spacer between GAPDH and a hypothetical protein (GAP2-IGS),

118    glutamine synthetase (GS), histone 3 (HIS3), nuclear rDNA internal transcribed

119    spacers (nrITS), mating type gene (MAT1-2-1), manganese-superoxide dismutase

120    (SOD2), and β-tubulin (TUB2).

121    It is known that the efficiency of PCR amplification and the distribution of

122  phylogenetic informativeness of a given marker varies among species complexes

123  (Hyde et al., 2013). Most studies on the utility and reliability of individual markers

124  come from the *Colletotrichum gloeosporioides* complex (Cai et al., 2009; Liu et al.,

125  2015; Sharma et al. 2013, Silva et al., 2012;2015; Vieira et al. 2017), and while the

126  *C. gloeosporioides* complex has been exhaustively studied in recent years,

127  recommendations on marker choice seem to be largely ignored. A classic case is the

128  utility of APN2/MAT-IGS, the most powerful marker to discriminate species within the

129  *C. gloeosporioides* complex: several recent studies described novel species within *C.*

130  *gloeosporioides* complex while excluding data from APN2/MAT-IGS (Costa et al.,

131  2018; Diao et al. 2017; Fu et al., 2019; Jayawardena et al., 2016; Oliveira et al.,

132  2018; Sharma et al., 2017; Silva et al., 2017; Sousa et al., 2018; Wang et al., 2019).

133  As mentioned above, this limits our ability to combine data from regional studies to

134  develop an accurate understanding of global diversity and phylogenetic relationships

135  within the genus.

136    While robust phylogenetic inference and reliable species delimitation relies on

137  the use of quality markers, studies on the performance of different molecular markers

138  for phylogenetic inference are missing for the majority of *Colletotrichum* species

139  complexes. In addition to the challenges discussed above, this also presents

140  practical problems for plant pathologists and ecologists who are looking to reliably

141  identify a large collection of isolates using molecular data. It is impractical for

142  researchers to sequence several loci, many of which may be of little phylogenetic

143  utility, across several hundred isolates simply for species identification. The aim of

144  the present study was to evaluate the phylogenetic informativeness of different

145  molecular markers used in *Colletotrichum* systematics and determine the optimal set

146     of markers for each species complex. From this, we hope to establish a consensus

147     on the minimal set of markers that can be used for *Colletotrichum* species

148     identification and delimitation and provide a practical reference for the large

149     community of researchers working on developing a better understanding of global

150     diversity, life history, and ecology of the genus.

151

152     **2. MATERIAL AND METHODS**

153     **2.1 Datasets**

154          We compiled several datasets to analyze the barcoding utility, phylogenetic

155     signal, and genealogical concordance of ACT, APN2, APN2/MAT-IGS, CAL, CHS-1,

156     GAPDH, GAP2-IGS, GS, HIS3, nrITS, and TUB2 for all *Colletotrichum* species

157     complexes described to date except for the *C. caudatum* complex, which was

158     excluded because nrITS was the only marker available for all species. These

159     datasets were compiled from published sequences retrieved from GenBank

160     (Supplementary File S1). Since testing the accuracy of prior species delimitations

161     were not the main focus of this study, we assumed that species boundaries

162     established in previous studies were accurate.

163          Thirteen species complexes were investigated in our study: *Colletotrichum*

164     *acutatum, C. boninense, C. dematium, C. destructivum, C. dracaenophilum, C.*

165     *gigasporum, C. gloeosporioides, C. graminicola, C. magnum, C. orbiculare, C.*

166     *orchidearum, C. spaethianum* and *C. truncatum*. Some species within each complex

167     were not included in the alignments due to the absence of sequences for several

168     markers since some of the analyses employed in the present work do not allow

169     missing data. Some markers were not analyzed due to a small number of species or

170     isolates with sequences available. The inclusion of these markers will drastically

171    reduce the number of species that can be included in each set of analyses (e.g.

172    GAPDH, HIS3, APN2 and APN2/MAT-IGS in the *C. graminicola* species complex).

173

174    **2.2 Multiple sequence alignment**

175          Multiple sequence alignments (MSA) of each locus were estimated individually

176    for each species complex. Sequences were compiled using the GenBank tool

177    implemented in MEGA 7 (Kumar et al., 2016). MSAs were estimated with the online

178    version of MAFFT 7 (Katoh et al., 2002; Katoh & Standley, 2013) using the G-INS-i

179    iterative refinement method and the 200PAM / k=2 nucleotide scoring matrix.

180    External gaps were trimmed in MEGA 7 before uploading MSAs to the GUIDANCE2

181    server (http://guidance.tau.ac.il/ver2/) (Sela et al., 2015) to access the alignment

182    confidence scores under the following parameters: MAFFT as the MSA algorithm;

183    max-iterate=0; pairwise alignment method=6mer; 100 bootstrap replicates. Unreliable

184    alignment regions were filtered by masking residues with scores below the lowest

185    cutoff, as proposed by Vieira et al. (2017). MSAs were converted to nexus format,

186    concatenated, and partitioned into a multilocus matrix using SequenceMatrix 1.8

187    (Vaidya et al., 2011). The number of invariable (I), variable (V), singletons (S) and

188    parsimony informative (PI) characters of the single locus alignments were calculated

189    using DnaSP 5.10 (Librado and Rozas,2009).

190

191    **2.3 DNA barcoding**

192          The effectiveness of markers to discriminate species within each species

193    complex was assessed by the barcode gap distance and intra/inter-specific distance

194    overlap (Hebert et al. 2003). Intra- and inter-specific distances were calculated for

195    each single locus alignment in MEGA 7. Single isolate species were removed from

196    alignments. Distances were calculated under the Kimura-2-parameter model,

197    allowing for substitution rates to differ among transitions and transversions, uniform

198    rates among sites, and gaps treated as pairwise deletions. Distance values were

199    sorted in Microsoft Excel Professional Plus 2016 and summary statistics were

200    calculated (maximum, minimum and mean distance). The barcode gap was

201    represented by the difference between the mean interspecific and intraspecific

202    distances (Hebert et al. 2003). The distance overlap percentage represents how

203    much the intraspecific distance overlaps with the interspecific distance and was

204    calculated as follows: max intraspecific distance ÷ max interspecific distance × 100.

205    Markers useful as barcodes will have a large barcode gap and a small intra/inter-

206    specific distance overlap.

207

208    **2.4 Assessment of phylogenetic informativeness**

209         The phylogenetic informativeness of markers commonly employed in

210    *Colletotrichum* systematics was estimated using the application PHYDESIGN

211    (Lopez-Giraldez and Townsend, 2011). Maximum likelihood (ML) trees were inferred

212    for each species complex using the concatenated alignments reduced to a single

213    representative isolate per species. Phylogenies were estimated in RAxML - HPC2

214    (Stawatakis, 2014) implemented on CIPRES Science Gateway portal

215    (https://www.phylo.org/portal2/home.action). ML tree searches were done assuming

216    the GTRGAMMA model and bootstrap support calculated with 1000 pseudoreplicates

217    (-m GTRGAMMA -p 12345 -k -f a -N 1000 -x 12345). ML trees were converted to

218    rooted ultrametric trees using the 'chronos' function in the ape package (Paradis et

219    al., 2004) using R Studio 1.1.442 (R Core Team, 2017). Trees were calibrated with

220    an arbitrary time scale with time=0 at the tips and time=1 at the root. The ultrametric

221    trees and the corresponding partitioned alignment were used as input files in

222    PHYDESIGN and the substitution rates were calculated using the program HyPhy

223    (Pond et al., 2005).  The substitution rates estimated by the maximum likelihood

224    algorithm used by HyPhy are nonsensical for some sites in the alignment resulting in

225    very recent 'phantom' peaks that have no biological meaning. These peaks are likely

226    the result of indels or ambiguous sites in the alignment, therefore these alignment

227    positions with poorly estimated substitution rates were excluded from some genes

228    prior to phylogenetic informativeness profiling at the recommendation of the authors

229    of PhyDesign (http://phydesign.townsend.yale.edu/faq.html). Since the markers

230    included in this study and in most systematic studies of *Colletotrichum* do not require

231    more than two sequencing reads to sequence the full length of the locus

232    (representing the same sequencing effort and cost), the phylogenetic informativeness

233    values (PIV) were calculated on a net basis. The variable PImax represents the time

234    in which a given marker reaches the maximum PIV and was used to determine the

235    divergence time in which the marker is most informative (Fong and Fujita, 2011). We

236    ranked the markers according to the PIV values and the usefulness of markers was

237    assessed through the profile shape: low and flat curves represent the least

238    informative markers; high and sharp peaks represent the most informative markers.

239        The percentage of markers that resolve a given species was estimated using a

240    Bayesian Concordance Analysis – BCA (Ané et al., 2007; Larget et al., 2010).

241    Although BCA is a coalescent-based method to estimate species trees (Ané et al.,

242    2007; Baum, 2007; Larget et al., 2010), this methodology can also be used to

243    quantify the proportion of markers that support a given clade, which is represented by

244    the concordance factor (CF). Individual locus trees were inferred in MrBayes 3.2.6

245    (Ronquist et al., 2012) implemented on the CIPRES cluster with four runs, each run

246    with four Markov chain Monte Carlo (MCMC) chains run for 10,000,000 generations,

247    sampled every 5,000 generations, totaling 2,001 trees per run. The frequency of

248    distinct topologies in the posterior distribution were summarized using mbsum

249    distributed with BUCKy 1.4.4 (Ané et al., 2007; Larget et al., 2010) skipping the first

250    25% of the trees as burn-in (-n 501). Primary concordance trees were estimated

251    using bucky and tree summary files output by mbsum were used as input files.

252    Concordance analyses were performed with a discordance factor (α) set at 1, four

253    MCMC chains, 1,000,000 generations, and the first 25% generations were discarded

254    as burn-in (-a 1 -k 4 -n 1000000 -c 4 -s1 23546 -s2 4564).

255         The Genealogical Sorting Index (GSI) was employed to identify the markers

256    that recover each species as monophyletic. GSI is an objective method that infers the

257    proportion of input trees for which a clade (species as applied here) are found to be

258    monophyletic and if the observed monophyly is greater than would be observed by

259    chance given the size of the data matrix (GSI=1 indicates monophyly) (Cummings et

260    al., 2008; Sakalidis et al., 2011), and can also be used to compare individual markers

261    according to their ability to discriminate species (Doyle et al., 2013). This

262    methodology can be applied to phylogenies inferred from a single locus as well as

263    multilocus analysis (Sakalidis et al., 2011). ML analyses were performed using the

264    single and multi-locus concatenated alignments. Species with a single isolate were

265    removed from the alignments. Analyses were carried out in RAxML as described

266    above with the number of pseudoreplicates reduced to 100 and outgroup isolates

267    specified prior to analysis. Rooted bootstrap trees were used as input files and each

268    tip was assigned to a species. GSIs were calculated using the GSI.py script and *P*-

269    value estimated from 100 permutations of each dataset (Cummings et al., 2008). The

270    GSI values were converted to heatmaps in the Heatmapper web server

271    (http://www.heatmapper.ca) to aid in visualization (Babicki et al., 2016).

272

273    **2.5 Selection of best minimal sets of markers**

274    The selection of a minimal set of optimal markers for robust phylogenetic

275    inference within each *Colletotrichum* species complex was based on results from the

276    phylogenetic informativeness profiling and GSIs. The markers were selected

277    according to the following criteria:

278    1) A minimum of three markers per complex based on ranking according to

279    PIV were selected. Three independent markers allow for the application of

280    the genealogical concordance phylogenetic species recognition criteria –

281    GCPSR (Dettman et al., 2003, Taylor et al. 2000), a commonly applied set

282    of criteria for phylogenetic species recognition in fungal systematics.

283    2) All species must be recognized as monophyletic by at least one of the

284    selected markers. GSIs were checked to confirm if each species in the

285    complex is recovered as monophyletic (GSI = 1) by at least one of the

286    selected markers.

287    3) ML trees were inferred from concatenated alignments of the three best

288    markers; if some species clade was poorly supported and/or all species

289    were not recovered as in the multilocus analyses with all markers

290    (unresolved relationships/polytomy), other markers were progressively

291    concatenated in decreasing order of phylogenetic informativeness until all

292    species were well resolved.

293

294       Once the best markers were chosen, the GSI was calculated for this

295   phylogeny to determine the level of species monophyly when only the best markers

296   are concatenated. The BCAs were also performed with this dataset to elucidate if the

297   species could be recognized by the majority of the selected markers.

298

299   **3. RESULTS AND DISCUSSION**

300   **3.1 Alignment statistics**

301       GAPDH, HIS3, and TUB2 were the most variable markers in the majority of

302   *Colletotrichum* species complexes, with PI characters ranging from 10—109, 11—82,

303   and 12—114, respectively (Table 1). PI characters for GS ranged from 63 to 93 and it

304   was the most variable marker within the *C. gigasporum* and *C. orbiculare* species

305   complexes. The APN2/MAT-IGS and GAP2-IGS, which are employed only in the *C.*

306   *gloeosporioides* species complex, had 192 and 115 PI characters, respectively, and

307   were the most variable markers within this complex. In contrast, nrITS presented the

308   fewest PI characters for most species complexes (0—36), followed by CHS-1 (3—45)

309   and ACT (3—63).

310       Most markers employed for *Colletotrichum* systematics comprise partial

311   sequences of orthologous protein-coding genes. These markers are composed of

312   introns flanked by long exons that are highly conserved. The GAPDH, GS, HIS3 and

313   TUB2 markers contain long intronic regions and APN2/MAT-IGS and GAP2-IGS

314   present long intergenic regions. This may explain why these markers are more

315   variable than the protein-coding loci. While protein-coding loci may be useful for

316   providing support along the backbone of the phylogeny within a species complex or

317   across the genus, markers with variable introns and intergenic sequences are

318   preferable for application at lower taxonomic levels (Schmitt et al., 2009).

319

320    **3.2 DNA barcoding feasibility**

321        GAPDH had the largest barcode gap distance in seven of the 11 species

322    complexes evaluated. The percentage overlap between intra- and interspecific

323    distances was less than 20%, with the exception of the *C. gigasporum* and the *C.*

324    *gloeosporioides* species complexes (28.6% and 29.2%, respectively). GS has the

325    highest barcode distance with the lowest overlap for the *C. gigasporum* (0.12 and

326    12.1%) and *C. orbiculare* (0.08 and 2.7%) complexes. APN2/MAT-IGS had the

327    largest barcode gap distance (0.15) and the smallest overlap percentage (3.26%)

328    within the *C. gloeosporioides* species complex, making it the best candidate barcode

329    locus for the complex HIS3 had the largest barcode gap (0.026) within the *C.*

330    *orchidearum* complex with a relatively low overlap (16.9%), although CHS-1 had a

331    comparable barcode gap (0.025) with only a slightly higher overlap (21.6%).  While

332    other markers, such as HIS3 and TUB2, are good candidate secondary barcode

333    markers in several complexes, nrITS was universally the poorest barcode candidate

334    with the lowest barcode gap distance within all species complexes.

335        Our results demonstrate that selecting a universal barcode marker for all

336    *Colletotrichum* species complexes among the markers currently being used is not

337    possible. An illustration of this is GAPDH. This marker is the best candidate barcode

338    marker for the majority of complexes, however it is among the worst barcode

339    candidates for the *C. gloeosporioides* species complex. While our results are not in

340    agreement with Cai et al. (2009), which carried out the first study evaluating the

341    markers to discriminate species within *C. gloeosporioides* species complex, only five

342    species were included in their analyses and the Musae and Kahawae clades (*sensu*

343    Weir et al. (2012)) were treated as single species. Moreover, while GAPDH was

344   chosen as the best marker relative to EF1α, ACT, CHS-1 and nrITS, the latter three

345   markers perform very poorly for species delimitation in the *C. gloeosporioides*

346   complex (Vieira et al., 2017) and, therefore, we did not include them in any of our

347   analyses. While GAPDH, with its large barcode gap and small overlap along with the

348   ease with which it can be amplified and sequenced, makes it among the best

349   barcode candidates across *Colletotrichum* as a whole, the selection of the best

350   barcode markers is dependent on the species complex.

351        The search for a universal barcode locus for the genus will require a

352   comparative analysis across the genomes of several species in all species

353   complexes. Intergenic sequences in syntenic regions of the genome are good

354   candidates if APN2/MAT-IGS and GAP-IGS are any indication. Intergenic sequences

355   may provide fast-evolving phylogenetic markers to be used for population genetic

356   and phylogenetic studies on fungal species complexes (Magain et al., 2017).

357

358   **3.3 Optimal markers for phylogenetic inference**

359        The most informative markers differ among the species complexes within

360   *Colletotrichum*. Net phylogenetic informativeness profiles and their respective

361   ultrametric trees are presented in Fig. 2. PIV and PImax values are summarized in

362   Table 1. The GSI analyses illustrate the ability of a given marker to recover species

363   as monophyletic is dependent on the species complex (Fig. 3, GSI values presented

364   in Supplementary File S2). Most species that were monophyletic in the multilocus

365   tree with all markers (GSI near to 1) were also recovered as monophyletic when only

366   the most informative markers were concatenated. In parallel, the BCAs revealed that

367   the proportion of markers supporting individual species-level clades (expressed as

368   concordance factors) increase when the less informative markers are removed from

369   the analyses.

370        Since the suite of markers differ by species complex and the performance of a

371   given marker differs among the species complexes, the results of the phylogenetic

372   informativeness profiling, GSI, Bayesian concordance analyses and selection of the

373   minimum set of markers for robust phylogenetic inference are summarized below for

374   each species complex.

375

376        **3.3.1 *Colletotrichum acutatum s. l.***

377        HIS3, GAPDH and TUB2 are the most phylogenetically informative markers

378   within the *C. acutatum* species complex, with PImax at 0.67, 0.78 and 0.99,

379   respectively (Fig. 2). All markers currently used for systematics of the C. *acutatum*

380   complex have an optimum inferential timescale varying from 0.67 to 0.99, which is

381   useful to resolve deeper relationships. However, the absence of markers with lower

382   PImax negatively impacts the our ability to infer relationships among recently

383   diverged species. This is evident by the low internal node support within Clade 1 and

384   Clade 2 *sensu* Damm et al. (2012a), which are the most recently diverged lineages

385   within the complex (Supplementary File S3).  Species are strongly supported in the

386   concatenated multilocus analysis (ML support ≥ 70%, Supplementary File S3),

387   whereas the relationship among them are poorly resolved even when only the best

388   markers were used to build the tree. In contrast, nodes throughout Clade 5 are

389   strongly supported. The inclusion of a marker with PImax near to 0.2 could improve

390   the support on the deeper nodes of the Clades 1 and 2.

391        Most species in *C. acutatum s. l.* could be resolved by the majority of the

392   markers only when the optimal markers are combined (Fig. 4B). Seven (*C. austral, C.*

393    *chrysanthemi*, *C. fioriniae*, *C. johnstonii*, *C. lupini*, *C. nymphaeae* and *C. tamarilloi*)

394    among the 17 multiple-isolate species were supported by the majority of individual

395    genes (CF≥0.66) in the BCA when all six loci were used in the analysis (Fig. 4A),

396    which means that four out of six markers support the monophyly of these species.

397    The CF increases when the analyses included only HIS3, GAPDH and TUB2, and

398    only four (*C. cosmi*, *C. costaricense*, *C. paranaense* and *C. phormii*) out of 17

399    multiple-isolate species presented CF<0.66 (Fig. 4B). The GSI (Fig. 3) indicates that

400    *C. melonis* is recovered by GAPDH and *C. phormii* by HIS3. The GSI for

401    *Colletotrichum costaricense* was nearly 1 for both GAPDH and TUB2 (0.95 and 0.99,

402    respectively) and less than 0.5 for all other markers. This is consistent with a CF of

403    0.63, indicating that nearly 2 of the 3 loci included in the analysis support the

404    monophyly of the species. *Colletotrichum paranaense* had a low GSI for most

405    markers and CF=0.22, which means that no individual marker fully supports this

406    species as monophyletic. However, *C. paranaense* is strongly supported in the

407    multilocus concatenated analysis and has a high GSI when analyses included all or

408    the best markers. These results clearly show that the multilocus trees masks the

409    incongruences among the individual gene trees or the use of markers with low

410    phylogenetic signal.

411        *Colletotrichum paranaense* was described by Bragança et al. (2016) as a

412    species closely related to *C. limetiicola* and *C. melonis* (Damm et al., 2012a).

413    Phylogenetic species recognition was not employed in the delimitation of *C.*

414    *paranaense*. The species was recognized based on the topology of the multilocus

415    tree and differed from *C. limetticola* and *C. melonis* by percentage identity according

416    to blastn searches. Our results suggest that the relationship among these species

417    and the delimitation of *C. paranaense* needs to be revisited.

418

419      **3.3.2 *Colletotrichum boninense s. l.***

420      CAL, TUB2 and GAPDH were the most informative markers to resolve species

421 within *C. boninense s. l.* (Fig. 2). All three markers are comparable in levels of

422 phylogenetic informativeness, with PIV ranging from 47 to 53. However, these

423 markers differ with respect to the optimal timescale at which they are informative.

424 CAL and TUB2 peaked at 0.58 and 0.63, respectively, providing more signal towards

425 the root of the phylogeny. In contrast, GAPDH reached maximum informativeness at

426 0.24, which is more useful for resolving species in this complex, since the divergence

427 epoch for most species is near 0.2.

428      Most species can be resolved by most of the markers according to the

429 concordance analyses that include all markers and only the best markers (CF>0.57

430 and >0.66, respectively) (Fig. 5B). GAPDH supports most species as monophyletic

431 (GSI 0.98—1) (Fig. 2), with the exception of three species (*C. catinaense*, *C.*

432 *constrictum* and *C. limonicola*). According to Damm et al. (2012b), all the species

433 within the *C. boninense* species complex can be identified by sequencing GAPDH. In

434 contrast, our results show that *C. catinaenese*, *C. constrictum* and *C. limonicola* are

435 weakly supported as monophyletic by GAPDH (GSI = 0.51, 0.69 and 0.59,

436 respectively) and could only be recovered as monophyletic in analyses of TUB2, HIS

437 or ACT. *Colletotrichum cymbidiicola* and *C. novae-zelandiae* were recovered as

438 monophyletic with GAPDH and *C. limonicola* by ACT and TUB2 (Fig. 3). Thus,

439 concordance factors for these species remained low even when the best markers

440 were used in the BCA. Additionally, *C. novae-zelandiae* was the only species for

441 which GSI<1 in both multilocus analyses (Fig. 3). *Colletotrichum catinaense* could be

442 resolved by two of the best markers (CAL and TUB) and by ACT, which led the CF to

443    shift from 0.56 in the BCA with all markers to 0.73 when only the best markers were

444    included (Fig. 5). The use of CAL, TUB2 and GAPDH leads to a slight increase in

445    support across the multilocus ML tree when compared with the tree reconstructed

446    with the whole dataset (Supplementary File S4), with the exception of *Colletotrichum*

447    *novae-zelandiae*, which remained poorly supported.

448

449    ### 3.3.3 *Colletotrichum dematium s. l.*

450    GAPDH, HIS3 and TUB2 presented high and sharp profiles (Fig. 2), with

451    PImax 0.05—0.06 and were the most useful markers to discriminate species within

452    the *C. dematium* species complex. HIS3 was the only marker able to discriminate all

453    species in the complex. nrITS, CHS-1 and ACT presented low and flat curves and

454    were the least informative markers (Fig. 2). *Colletotrichum anthrisci* and *C. circinans*

455    could be resolved by all markers (CF=1, Fig. 6), while *Colletotrichum spinaceae* was

456    resolved as monophyletic by all markers except nrITS. The exclusion of nrITS in the

457    BCA increased the CF from 0.86 to 1. The relationship between *Colletotrichum*

458    *dematium* and *C. lineola* was not clearly resolved in the multilocus analyses.

459    Although *C. dematium* isolates were well supported as monophyletic in the tree

460    inferred from all markers, *C. lineola* isolates remained paraphyletic (Supplementary

461    File S5). When only GAPDH, HIS3 and TUB2 were considered, the isolates of *C.*

462    *dematium* and *C. lineola* were placed together in a poorly supported clade.

463    *Colletotrichum lineola* was recovered as monophyletic only by HIS3 (GSI=1). Both

464    ACT and GAPDH could also resolve this species, albeit with low GSI (0.83 and 0.73

465    respectively). In contrast, these three markers recover *C. dematium* with high level of

466    monophyly (GSI 0.97—1). *Colletotrichum lineola* presented low monophyly level in all

467  multilocus GSI analysis and *C. dematium* was not monophyletic when only the best

468  markers were considered.

469  In the Damm et al. (2009) study, the isolates of *C. dematium* and *C. lineola* are

470  separated into two short-branch clades. They chose to retain both taxa, but raised

471  the hypothesis that *C. dematium* and *C. lineola* are different populations within the

472  same species. We performed the GSI analyses with the multilocus trees combining

473  *C. dematium* and *C. lineola* in the same group (data not shown). The group was

474  recovered as monophyletic when all markers and only the best markers are

475  concatenated (GSI=1 and 0.99 respectively), which support the hypothesis that *C.*

476  *dematium* and *C. lineola* are the same species. *Colletotrichum eryngiicola, C.*

477  *hemerocallidis, C. insertae, C. quiquefoliae* and *C. sonchicola* were placed together

478  with *C. dematium* and *C. lineola* in a polytomous clade in Samarakoon et al. (2018)

479  study, which indicates that those species may also be members of the group *C.*

480  *dematium/ C. lineola*. These species were not included in our analyses due to

481  absence of HIS3 and TUB2 sequences. In the future, species boundaries within this

482  lineage need to be revisited.

483

484  **3.3.4 *Colletotrichum destructivum* s. l.**

485  HIS3, TUB2 and GAPDH were the most informative markers to resolve

486  species within *C. destructivum* complex (Fig. 2). GAPDH and HIS3 possess the

487  phylogenetic signal to resolve shallow clades (PImax=0.28 and 0.37 respectively),

488  whereas TUB2 performs better in the deep branches (PImax=0.73). HIS3 was the

489  most informative marker and was able to six (*C. destructivum, C. lentis, C. lini, C.*

490  *tabacum, C. utrechtense* and *C. vignae*) out of the nine species analyzed as

491  monophyletic (GSI 0.98—1) (Fig. 3). Although GAPDH is one of the best markers,

492    only four species (*C. fuscum*, *C. lentis*, *C. tabacum* and *C. vignae*) were highly

493    monophyletic in the topologies provided by this locus (GSI 0.9—1) (Fig. 3).

494    *Colletotrichum americae-borealis* was not monophyletic in the topology of any single

495    or multi-locus trees (GSI 0.34—0.89). Four (*C. destructivum*, *C. lentis*, *C. tabacum*,

496    *C. utrechtense*) out of nine multiple-isolates species could be recovered by most of

497    the genes in the concordance tree (CF≥0.66) inferred from all markers (Fig. 7A).

498    These same four species plus *C. vignae* were recovered by most of the genes

499    (CF≥0.66) when only the optimal markers were considered. Moreover, the CF and

500    the bootstrap supports of the internal nodes increased when only the best markers

501    were combined (Fig. 7 and Supplementary File S6).

502         The clade comprised by *C. americae-borealis* and *C. lini* forms a polytomy

503    when only the three best markers were used (Supplementary File S6). We also

504    tested the inclusion of ACT or nrITS in the concatenated dataset. However, this clade

505    remains a polytomy and can only be resolved when all markers are concatenated.

506    The overall GSI values also reduced when ACT or nrITS was included in the

507    multilocus analysis (data not shown). Thus, it is not reasonable to include ACT and

508    nrITS in the analysis due to their low phylogenetic informativeness. Other markers

509    with greater phylogenetic signal need to be tested for this complex in order to resolve

510    the relationships among the unresolved clades and clarify the species identities.

511

512    **3.3.5 *Colletotrichum dracaenophilum s. l.***

513         TUB2, GAPDH and HIS3 presented the highest PIV (58, 52 and 36

514    respectively) and were the most informative markers to distinguish species within the

515    *C. dracaenophilum* complex (Fig. 2). All the markers peaked around the same

516    timescale (PImax 0.39—0.46) and provide robust support for relationships in both

517   deep and shallow nodes. Although TUB2, GAPDH and HIS3 were the most

518   informative markers, all the markers are informative enough to discriminate species

519   in *C. dracaenophilum s. l.* (GSI 0.99—1) (Fig. 3). This is corroborated by the high

520   concordance factors (CF=1) in BCAs (Fig. 8) and 100% support in multilocus

521   analyses (Supplementary File S7) when all or the best markers were combined.

522

523       **3.3.6 *Colletotrichum gigasporum s. l.***

524       GS, CAL and GAPDH were the most powerful markers to discriminate species

525   within *C. gigasporum s. l.* (Fig. 2), with PImax at 0.1, 0.44 and 0.28 respectively.

526   Although GS, CAL and GAPDH were the most informative markers, all the markers

527   and the concatenated datasets were able to discriminate all species within the *C.*

528   *gigasporum* complex (GSI 0.99—1) (Fig. 3), as reported by Liu et al. (2014). All

529   species presented maximum CF (Fig. 9) and ML support (Supplmentary File S8)

530   independently of which set of markers was included in these analyses.

531

532       **3.3.7 *Colletotrichum gloeosporioides s. l.***

533       APN2/MAT-IGS, GAP2-IGS and APN2 were the most informative markers to

534   separate species within the *C. gloeosporioides* species complex (Fig. 2). The

535   informativeness profiles indicate APN2/MAT-IGS and APN2 are of peak

536   informativeness at 0.19 and 0.17, respectively, and are informative for shallow

537   divergences, whereas GAP2-IGS has an optimal timescale at 0.42 and provides

538   more signal for resolving deep nodes. This range in values for peak informativeness

539   led to high support at both deep and shallow nodes when only the most informative

540   markers were used for phylogenetic inference (Supplementaray File S9). APN2/MAT-

541   IGS could separate all the species included in the analyses, although only *C.*

542    *fragariae*, *C. queenslandicum*, *C. siamense* and *C. tropicale* reached a GSI of 1 (Fig.

543    3). Moreover, all species were recovered as monophyletic in all multilocus analyses.

544    *Colletotrichum siamense* was the only species that could not be recovered by most of

545    the markers, and its CF reduced from 0.59 in all-markers analysis to 0.38 in the best-

546    markers analysis (Fig. 10B). The CF reduced due to removing CAL and TUB2, which

547    also recover *C. siamense* as monophyletic (GSI=0.99).

548          ACT, CHS-1 and nrITS were not included in our analyses because these

549    markers were previously reported as the worst markers to distinguish species in *C.*

550    *gloeosporioides* complex (Vieira et al., 2017). In the present study, rather than

551    restating the results of Vieira et al. (2017), we evaluated if the seven markers

552    proposed by those authors are needed for diversity studies and species assignment.

553    We determined sequences of APN2/MAT-IGS, GAP2-IGS and APN2 are sufficient to

554    resolve the species in *C. gloeosporioides s. l.*. However, sequences of GAP2-IGS

555    and APN2 are not available for all species within *C. gloeosporioides s.l..* These data

556    should be generated for all species within the species complex.

557          *Colletotrichum gloeosporioides s. l.* is the most studied species complex in the

558    genus. More than 10 different markers have been used for systematics and

559    taxonomy of the *C. gloeosporioides* complex over the last 10 years. However, there

560    is no agreement about which markers should be used for species recognition. The

561    first attempt to determine the best markers was done by Cai et al. (2009), in which

562    GAPDH was the best marker to discriminate species in *C. gloeosporioides s. l.*, and

563    the set composed by ACT, GAPDH, nrITS and TUB2 was recommended for

564    multilocus analysis. Weir et al. (2012) highlighted that although GAPDH is one of the

565    most effective markers to distinguish species within *C. gloeosporioides* complex, the

566    combination with GS is necessary to distinguish some species. Our study shows that

567    GAPDH was the least informative marker among the ones included in this study for

568    *C. gloeosporioides s. l.*. This marker was the least variable (Table 1), and had the

569    smallest barcode gap and the largest overlap distance (Fig. 1). *Colletotrichum*

570    *fragariae* was the only species recovered as monophyletic (GSI=1) by this marker

571    (Fig 2). Based on this, GAPDH is considered one of the worst barcode candidates for

572    the *C. gloeosporioides* complex among the markers tested in the present study.

573        More recently, several studies demonstrate the singular ability of APN2/MAT-

574    IGS to discriminate species in the *C. gloeosporioides* complex (Sharma et al., 2013,

575    2015; Vieira et al., 2014), although others previously demonstrated the utility of this

576    marker (Du et al., 2005; Rojas et al., 2012; Silva et al., 2012). The main issue in

577    using APN2/MAT-IGS was the splitting out of *C. siamense* in a species complex, in

578    which several monophyletic lineages could be revealed by the phylogeny inferred by

579    this marker (Sharma et al., 2013). Although the multilocus analysis had been done,

580    the lineages identity was confirmed only based on the APN2/MAT-IGS, since this

581    marker performs good as well as the multilocus matrix. Other studies use the same

582    criteria to discriminate species and several species within *C. siamense s. l.* were

583    described (Sharma et al., 2015; Vieira et al., 2014). Later, Liu et al. (2016) use

584    coalescent methods for phylogenetic species delimitation and synonymize all those

585    species in the complex into *C. siamense*. It was revealed incongruences among the

586    APN2/MAT-IGS tree and other individual gene trees. The study clarifies how

587    multilocus analyses can mask discordances among individual gene trees and lead to

588    species misidentification.

589        The combination of APN2/MAT-IGS and GS were proposed as the barcode to

590    delimit species within *C. gloeosporioides* complex (Liu et al., 2015). These two

591    markers were collectively powerful enough to discriminate the 22 species included in

592    the study and produced the same topology as that inferred from a 6 marker

593    multilocus dataset (ACT, CAL, GAPDH, GS, nrITS and TUB2). We tested the

594    combination of APN2/MAT-IGS with each remaining marker individually and all

595    resulting trees were similar in topology (data not shown). Some species, such as *C.*

596    *siamense*, are polyphyletic in our GAPDH and GS trees, which belies the

597    incongruence with the multilocus and the other individual gene trees as presented by

598    Silva et al. (2012) and Vieira et al. (2017).  Mating-type associated markers, such as

599    APN2/MAT-IGS and MAT1-2, had fast evolutionary rates and high variability, and

600    may dominate the topology of multilocus trees (Liu et al., 2016). Since the

601    combination of APN2/MAT-IGS with any other marker produces similar topologies,

602    the combination with any other marker besides GAPDH and GS is preferable in order

603    to avoid inconsistencies in species delimitation within the *C. gloeosporioides*

604    complex. While we find that APN2 is another informative marker that should be

605    incorporated into diversity studies of the *C. gloeosporioides* complex, the proximity of

606    this marker to the APN2/MAT-IGS region suggests it is part of the same linkage

607    group and thus may not represent an independent sample of an organisms

608    evolutionary history. If this is the case, linking substitution models and inferred trees

609    across these two loci would be required and the addition of another locus may be

610    necessary to fulfill the expectations under GCPSR.

611          Although the majority of *Colletotrichum* species belong to the *C.*

612    *gloeosporioides* complex, identification and description of taxa within this complex is

613    the most problematic. Sequences of the best markers, mainly APN2 and GAP2-IGS,

614    have not been sequenced for the majority of species, which prohibits the detection of

615    these species using these markers. A priority should be to generate sequences of the

616    best markers for all species to avoid further misidentification and the introduction of

617    dubious taxa. Until this happens, the description of new species within the *C.*

618    *gloeosporioides* complex will likely require sequencing multiple additional markers to

619    be certain of their novelty.

620

621         **3.3.8 *Colletotrichum graminicola s. l.***

622         TUB2 was the most informative marker, followed by nrITS and ACT (Fig. 2). In

623    contrast, CHS-1 was ranked as the worst marker for the *C. graminicola* species

624    complex. GSI and BCA analyses could not be performed for this complex because

625    only one isolate per species could be included in the analyses. Other markers such

626    as APN2, APN2/MAT-IGS, GAPDH, HIS3, MAT1-2 and SOD2 were also used to

627    identify species within *C. graminicola* complex (Cannon et al., 2012; Crouch et al.

628    2009b, c; Crouch and Tomaso-Peterson, 2012; Du et al., 2005; Moriwaki and

629    Tsukiboshi, 2009; O'Connell et al., 2012; Tao et al., 2013). However, sequences for

630    these markers are missing for several species. According to our results for other

631    complexes, APN2, APN2/MAT-IGS, GAPDH and HIS3 are likely to be more powerful

632    markers than those that we could include in the analyses for the *C. graminicola*

633    complex. Sequences for these markers need to be generated for isolates of several

634    species in order to establish a better set of optimal markers for phylogenetic

635    inference and species discrimination in the *C. graminicola* complex.

636

637         **3.3.9 *Colletotrichum magnum s. l.***

638         HIS3, TUB2 and GAPDH were the most informative markers to discriminate

639    species within the *C. magnum* complex (Fig. 2). Although HIS3 presented the highest

640    PIV (26), TUB2 presented a higher PImax (0.52 versus 0.44) and is able to separate

641    more species than HIS3. Only *C. brevisporum* and *C. magnum* were accounted for

642    by more than one isolate and could be analyzed in the BCA and GSI analyses (Fig.

643    3). *Colletotrichum brevisporum* reach high levels of monophyly only in analyses of

644    the concatenated dataset (GSI 0.98—1, Fig. 3) and only one gene can fully resolve

645    this clade (CF=0.4, Fig 11). In contrast, *C. magnum* is recovered as monophyletic by

646    the multilocus dataset and by the best markers. However, the CF reached 1 only

647    when the best markers were concatenated (Fig 11). *Colletotrichum magnum*

648    remained highly supported when only the best markers were used to reconstruct the

649    phylogeny of the *C. magnum* complex (Supplementary File S10B) with only a slight

650    decrease in support at some internal nodes. The species *C. liaoningense* was not

651    included in the analyses due to sequence deposition errors detected by Damm et al.

652    (2019) from the study were *C. liaoningense* was described (Diao et al., 2017). Damm

653    et al. (2019) concluded that *C. liaoningense* needs to be revisited.

654

655          **3.3.10 *Colletotrichum orbiculare* s. l.**

656          GS stood out from all markers as the most informative (PIV=69) in the *C.*

657    *orbiculare* complex, followed by HIS3 and GAPDH (PIV=37 and 32 respectively) (Fig.

658    2). GS peaks at 0.53 and can discriminate the majority of species and provide robust

659    support for the relationships among them, since the divergence time for most species

660    is about 0.3. On the other hand, HIS3 and GAPDH are useful to discriminate some

661    recently diverged species (PImax=0.13 and 0.17 respectively) such as *C. orbiculare*

662    × *C. sidae* and *C. trifolii* × *C. malvarum*. GS was the marker that recovered more

663    species as monophyletic (GSI=1), with the exception of *C. sidae* (GSI=0.89) (Fig. 3).

664    *Colletotrichum lindemuthianum* was the only species that was not recovered as

665    monophyletic when all or the best markers were concatenated (both GSI=0.78),

666    which is likely due to the variability within this species as currently circumscribed.

667 Damm et al. (2013) split *C. lindemutianum* into two different lineages (*C.*

668 *lindemuthianum* 1 and 2), which are observed in the nrITS and GS trees. In our

669 multilocus trees, the isolates CBS133.57 and CBS131.57 where moved to the clade

670 *C. lindemuthianum* 1 and CBS150.28 to the clade *C. lindemuthianum* 2. This result is

671 discordant with that in Damm et al. (2013) and we conclude that the terminal clades

672 within *C. lindemuthianum* represent intraspecific variability. Thus, both *C.*

673 *lindemuthianum* lineages were not considered separate species in our analysis. The

674 CF of both deep and shallow nodes increases significantly when only the best

675 markers were analyzed (Fig.12).  Reducing the set of markers does not cause

676 significant differences in the tree topology or clade support (Supplementary file S10).

677

678 **3.3.11 *Colletotrichum orchidearum s. l.***

679 TUB2 and HIS3 were the most powerful genes capable of discriminating

680 species within the *C. orchidearum* complex (PIV=17 and 16 respectively), followed by

681 CHS-1 (PIV=12) (Fig. 2). TUB2 and HIS3 are good markers to discriminate all

682 species in the complex and support the relationships among them (PImax=0.8 and

683 0.67 respectively), whereas CHS-1 can help to distinguish and support recently

684 diverged species (PImax=0.47). All species were recovered with high levels of

685 monophyly with TUB2 (GSI=0.99—1), and also by both multilocus analyses

686 (GSI=0.97—1) (Fig. 3). *Colletotrichum musicola* was the only species that could be

687 recovered by all markers evaluated (CF=1) (Fig. 13). The CFs increased when the

688 less informative markers were removed from the analysis and most species were

689 recovered by all markers (CF=1), with the exception of *C. cliviicola and C. plurivorum*

690 (CF≥0.66). Damm et al. (2019) also reported that some clades were not supported by

691 some of the individual gene analyses. All deep and shallow nodes retained high

692    support when only the best markers were combined in our analyses (Supplementary

693    File S12B).

694

695    **3.3.12 *Colletotrichum spaethianum s. l.***

696    GAPDH, TUB2 and ACT were the best markers to separate species within the

697    *C. spaethianum* complex (Fig. 2). These markers peak at approximately 0.4 (PImax

698    0.47—0.77), which is in the range of where most species in the complex diverge

699    indicating they can separate the majority of species in this complex. However,

700    GADPH is the only marker that recovers all species as monophyletic (GSI 0.98—1)

701    (Fig. 3). All species reached complete monophyly (GSI=1) in both multilocus trees.

702    *Colletotrichum liriopes* and *C. spaethianum* could be recovered by most of the genes

703    only when the best markers are considered in the BCA (CF=0.77 and 0.70,

704    respectively) (Fig. 14B). All species remain strongly supported when only the best

705    markers were concatenated (Supplementary Figure S13B).

706    HIS3 appears to be a good marker in several other *Colletotrichum* complexes

707    according to the present study, and may be among the three best markers for *C.*

708    *spaethianum* complex. However, sequences of HIS3 and other markers such as

709    CHS-1 and CAL are available only for one isolate of several species and we cannot

710    include them in our analyses.  Sequences of these markers need to be generated for

711    other isolates to determine if one of these markers could be used to substitute for

712    ACT, which does not perform well for separating *Colletotrichum* species.

713

714    **3.3.13 *Colletotrichum truncatum s. l.***

715    GAPDH was clearly the most informative marker (PIV=57), followed by TUB2

716    and ACT (PIV=33 and 23 respectively) (Fig. 2). All markers peaked above 0.4 and

717    were able to discriminate most species. *Colletotrichum acidae* and *C. curcumae* were

718    monophyletic (GSI=1) in all single and multilocus datasets (Fig. 3) with CFs equal to

719    1 in analyses of both all and best markers BCA (Fig. 15). These species were also

720    supported by maximum bootstrap values in the multilocus trees (Supplementary File

721    S14). *Colletotrichum corchorum-capsularis* and *C. truncatum* were not recovered as

722    monophyletic by any dataset. Moreover, *C. truncatum* is paraphyletic or polyphyletic

723    in the multilocus trees (Fig. 15, Supplementary File S14), which leads us to the

724    conclusion that *C. corchorum-capsularis*, as circumscribed by Niu et al. (2016),

725    cannot be recognized as a species distinct from *C. truncatum*. Our results strongly

726    suggest that *C. corchorum-capsularis* and *C. truncatum* may be the same species

727    and isolates of both species were placed together in a clade with high CF (≥0.96) in

728    the BCAs and maximum support in the multilocus analyses (Fig. 15, Supplementary

729    File S14). Additional work is needed using the best markers for the *C. truncatum*

730    complex and objective species recognition methods to determine the taxonomic

731    status species boundaries of *C. corchorum-capsularis* and *C. truncatum*.

732

733    **4. Conclusions**

734        We used phylogenetic informativeness profiling, maximum likelihood and

735    coalescent-based phylogenetic analyses, measures of barcode utility, and

736    genealogical sorting indices to assess the performance of the several molecular

737    markers used in *Colletotrichum* systematics and taxonomy across all known species

738    complexes. While HIS3, GAPDH, and TUB2 were among the best markers for most

739    of the complexes, the optimal set of markers is not always the same across all

740    complexes. ACT, CHS-1 and nrITS were the worst markers and, as previously

741    proposed for the *C. gloeosporioides* complex (Vieira et al., 2017), they can be

742     discarded from the phylogenetic analysis of almost all species complexes. ACT was

743     retained in the set of best markers for the *C. graminicola*, *C. spaethianum* and *C.*

744     *truncatum* complexes to achieve a minimum of three markers as proposed in the

745     methodology of the present study. However, few markers were included for these

746     complexes due to missing data, therefore additional markers need to be sequenced

747     and their performance evaluated. The analyses of *C. caudatum, C. dematium, C.*

748     *graminicola* and *C. spaethianum* complexes were the most impacted by the

749     excessive amount of missing data for the majority of the markers, which highlights

750     the importance of selecting a standard set of markers to delimit species. Similarly,

751     several isolates and/or species were not included in the marker analyses for *C.*

752     *gloeosporioides s.l.* due to selective data acquisition by different research groups. It

753     is not clear how the inclusion of sequences from these isolates might impact our

754     results. Sequences need to be generated for these markers and/or species to

755     provide more decisive results. We have also identified species complexes, which will

756     need to be revisited in the future, in which it appears species have been misidentified

757     (*C. acutatum*, *C. dematium* and *C. truncatum)*.

758          Selecting the optimal markers to sequence for biodiversity studies on

759     *Colletotrichum* will impact *Colletotrichum* studies in a few ways. First, species

760     recognition will likely be more accurate and robust by avoiding the confounding effect

761     of including markers with low phylogenetic signal in the analyses. Secondly,

762     phylogenetic studies of *Colletotrichum* will become more economical, since

763     sequencing markers with low informativeness represents a low return on investment.

764     Finally, if research groups take guidance from this study, we are more likely to see a

765     consensus developed on the data acquired for phylogenetic studies on *Colleotrichum*

766     and we will be closer to a global assessment with combinable data.

767   Researchers around the world continue to have an interest in documenting the

768  diversity of *Colletotrichum* species associated with economically important plant

769  species. However, this work is labor intensive and expensive because several

770  hundred isolates are typically screened and a paucity of distinctive morphological

771  characters necessitates DNA sequencing. The expense has been unnecessarily

772  compounded by the lack of an objective and comprehensive assessment of the utility

773  of existing markers for phylogenetic inference and species identification/delimitation

774  and the lack of a consensus on the markers to be used. We hope the results

775  presented here will help to address this problem. While the optimal markers differ by

776  species complex, our results provide some guidance on the most efficient path to

777  document and describe diversity within *Colletotrichum*. Our results also show that for

778  the accurate identification and delimitation of *Colletotrichum* species, a small set of

779  markers with strong phylogenetic signal is more suitable than a large set including

780  markers with both weak and strong phylogenetic signal. GAPDH is among the

781  optimal set of markers for 10 of the 13 species complexes in *Colletotrichum*, followed

782  by TUB2 (10 of 13), and HIS3 (7 of 13). Therefore, GAPDH is a good marker to

783  sequence for initial diversity screening and assigning isolates to a species complex

784  because data for this marker is available for the majority of species within the genus.

785  However, selection of additional markers for phylogenetic inference and species

786  delimitation will depend on the species complex.

787   Finally, while we recommend the optimal markers for species recognition

788  within *Colletotrichum* in order to improve diversity studies in the genus, our

789  understanding of evolutionary relationships among species remains poorly resolved.

790  Improving our understanding of relationships among taxa within *Colletotrichum* will

791  require more robust genomic sampling. Genome sequencing is underway for many

792    species of *Colletotrichum*, however a comprehensive phylogenomic study of the

793    genus is needed.

794

**Acknowledgments**

807

**Glossary**

**Appressorium**: specialized cell produced by some phytopathogenic fungi which is used to infect plant hosts.

**Conidium**: asexual spore of Ascomycota and Basidiomycota.

**Endophytic fungi**: fungi that grow inside the plant tissues without causing disease symptoms.

**Phytopathogenic fungi**: fungi that cause plant diseases.

***Sensu lato* (*s. l.*)**: taxonomic terminology used to reference species complexes (*C. acutatum* species complex = *C. acutatum s. l.*).

817     ***Sensu stricto* (*s. s.*)**: when is necessary to refer the species with the same name of

818     the complex (*C. acutatum s. s.* is a species within *C. acutatum s. l.*).

819     **Species complex**: major clades strongly supported within *Colletotrichum* genus tree.

820     These clades include phylogenetic species closely related which most are

821     indistinguishable based on phenotypical characters (e.g. conidial and appressorial

822     shape and size, growth rate, color of colonies). Species complexes get the same

823     name of the species within them that is more known or that was firstly described.

824     In some cases, members within a given species complex share peculiar conidial

825     characteristics: *C. acutatum* – conidia with acute ends; *C. boninense* – presence of a

826     prominent scar (hilum) at the base of the conidium; *C. caudatum* – conidia with a

827     filiform appendage at the apex; *C. gigasporum* – longest and widest conidia within

828     the genus.

829

830     **References**

831     Ané, C., Larget, B., Baum, D.A., Smith, S.D., Rokas, A., 2007. Bayesian estimation

832       of concordance among gene trees. Molecular Biology and Evolution 24, 412–

833       426.

834     Babicki, S., Arndt, D., Marcu, A., Liang, Y., Grant, J.R., Maciejewski, A., Wishart,

835       D.S., 2016. Heatmapper: web-enabled heat mapping for all. Nucleic Acids

836       Research 44, W147–W153.

837     Baum, D.A., 2007. Concordance trees, concordance factors, and the exploration of

838       reticulate genealogy. Taxon 56, 417–426.

839     Bragança, C.A.D., Damm, U., Baroncelli, R., Massola Júnior, N.S., Crous, P.W.,

840       2016. Species of the *Colletotrichum acutatum* complex associated with

841       anthracnose diseases of fruit in Brazil. Fungal Biology 120, 547–561.

842    Cai, L., Hyde, K.D., Taylor, P.W.J., Weir, B.S., Waller, J., Abang, M.M., Zhang, J.Z.,

843        Yang, Y.L., Phoulivong, S., Liu, Z.Y., Prihastuti, H., Shivas, R.G., McKenzie,

844        E.H.C., Johnston, P.R., 2009. A polyphasic approach for studying *Colletotrichum*.

845        Fungal Diversity 39, 183–124.

846    Cannon, P.F., Damm, U., Johnston, P., Weir, B.S., 2012. *Colletotrichum* – Current

847        status and future directions. Studies in Mycology 73, 181–213.

848    Cao, X.; Xu, X.; Che, H.; West, J.S.; Luo, D., 2019. Three Colletotrichum Species,

849        Including a New Species, are Associated to Leaf Anthracnose of Rubber Tree in

850        Hainan, China. Plant Disease 103, 117–124.

851    Corda, A.C.I., 1831. Die Pilze Deutschlands. In: Deutschlands Flora in Abbildungen

852        nach der Natur mit Beschreibungen (Sturm, J, ed.). Sturm, Nürnberg vol. 3, Abt.

853        12: 33–64, tab. 21–32.

854    Costa, J.F.O., Kamei, S.H., Silva, J.R.A., Miranada, A.R.G.S., Netto, M.B., Silva,

855        S.J.C., Correia, K.C., Lima, G.S.A., Assunção, I.P., 2018. Species diversity of

856        *Colletotrichum* infecting *Annona* spp. in Brazil. European Journal of Plant

857        Pathology. https://doi.org/10.1007/s10658-018-01630-w.

858    Crouch, J.A., Clarke, B.B., Hillman, B.I, 2009a. What is the value of ITS sequence

859        data in *Colletotrichum* systematics and species diagnosis? A case study using

860        the falcate-spored graminicolous *Colletotrichum* group. Mycologia 101, 648–656.

861    Crouch, J.A., Beirn, L.A., Cortese, L.M., Bonos, S.A., Clarke, B.B., 2009b.

862        Anthracnose disease of switchgrass caused by the novel fungal species

863        *Colletotrichum navitas.* Mycological Research 113, 1411–1421

864    Crouch, J.A., Tredway, L.P., Clarke, B.B. and Hillman, B.I., 2009c. Phylogenetic and

865        population genetic divergence correspond with habitat for the pathogen

866     *Colletotrichum cereale* and allied taxa across diverse grass communities.

867     Molecular Ecology 18, 123–135.

868     Crouch, J.A., Tomaso-Peterson, M., 2012. Anthracnose disease of centipedegrass

869     turf caused by *Colletotrichum eremochloae*, a new fungal species closely related

870     to *Colletotrichum sublineola*. Mycologia 104, 1085–1096.

871     Cummings, M.P., Neel, M.C., Shaw, K.L., 2008. A genealogical approach to

872     quantifying lineage divergence. Evolution 62, 2411–2422.

873     Damm, U., Woudenberg, J.H.C., Cannon, P.F., Crous, P.W., 2009. *Colletotrichum*

874     species with curved conidia from herbaceous hosts. Fungal Diversity 39, 45–87.

875     Damm, U., Cannon, P.F., Liu, F., Barreto, R.W., Guatimosim, E, Crous,P.W., 2013.

876     The *Colletotrichum orbiculare* species complex: Important pathogens of field

877     crops and weeds. Fungal Diversity 61, 29–59.

878     Damm, U., Cannon, P.F., Woudenberg, J.H.C., Crous, P., 2012a. The *Colletotrichum*

879     *acutatum* species complex. Studies in Mycology 73, 37–113.

880     Damm, U., Cannon, P.F., Woudenberg, J.H.C., Johnston, P.R., Weir, B.S., Tan, Y.P.,

881     Shivas, R.G., Crous, P.W., 2012b. The *Colletotrichum boninense* species

882     complex. Studies in Mycology 73:1–36.

883     Damm, U., Sato, T., Alizadeh, A., Groenewald, J.Z., Crous,P.W., 2019. The

884     *Colletotrichum dracaenophilum*, *C. magnum* and *C. orchidearum* species

885     complexes. Studies in Mycology 92, 1–46.

886     Dean, R., Van Kan, J.A.L., Pretorius, Z.A., Hammond-Kosack, K.E., Di Pietro, A.,

887     Spanu, P.D., Rudd, J.J., Dickman, M., Kahmann, R., Ellis, J., Foster, G.D., 2012.

888     The Top 10 fungal pathogens in molecular plant pathology. Molecular Plant

889     Pathology 13, 414–430.

890    Dettman, J.R., Jacobson, D.J., Taylor, J.W., 2003. A multilocus genealogical

891        approach to phylogenetic species recognition in the model eukaryote

892        *Neurospora*. Evolution 57, 2703–2720.

893    Diao, Y.-Z., Zhang, C., Liu, F., Wang, W.-Z., Liu, L., Cai, L., Liu, X.-L., 2017.

894        *Colletotrichum* species causing anthracnose disease of chili in China. Persoonia

895        38, 20–37.

896    Doyle, V.P., Oudemans, P.V., Rehner, S.A., Litt, A., 2013. Habitat and host indicate

897        lineage identity in *Colletotrichum gloeosporioides s. l.* from wild and agricultural

898        landscapes in North America. PLoS ONE 8, e62394.

899    Du, M., Schardl, C., Nuckles, E., Vaillancourt, L., 2005. Using mating-type gene

900        sequences for improved phylogenetic resolution of *Collectotrichum* species

901        complexes. Mycologia 97, 641–658.

902    Fong, J.J., Fujita, M.K., 2011. Evaluating phylogenetic informativeness and data-type

903        usage for new protein-coding genes across Vertebrata. Molecular Phylogenetics

904        and Evolution 61, 300–307.

905    Fu, M., Crous, P.W., Bai, Q., Zhang, P.F., Xiang, J., Guo, Y.S., Zhao, F.F., Yang,

906        M.M., Hong, N., Xu, W.X. and Wang, G.P., 2019. *Colletotrichum* species

907        associated with anthracnose of *Pyrus* spp. in China. Persoonia 42, 1–35.

908    Hebert, P.D., Cywinska, A., Ball, S.L., 2003.  Biological identifications through DNA

909        barcodes. Proceedings of the Royal Society of London Series B 270, 313–321.

910    Hyde KD, Cai L, McKenzie EHC, Yang YL, Zhang JZ, Prihastuti, H., 2009.

911        *Colletotrichum*: a catalogue of confusion. Fungal Diversity 39, 1–17.

912    Hyde, K.D., Udayanga, D., Manamgoda, D.S., Tedersoo, L., Larsson, E., Abarenkov

913        K., Bertrand, Y.J.K., Oxelman, B., Hartmann, M., Kauserud, H., Ryberg, M.,

914        Kristiansson, E., Nilsson, R.H., 2013. Incorporating molecular data in fungal

915    systematics: A guide for aspiring researchers. Current Research in

916    Environmental and Applied Mycology 3, 1–32.

917  Jayawardena, R.S., Yan, J., Hyde, K.D., Zhang, G.,2016. Morphological and

918    molecular characterization of *Colletotrichum* species of strawberry in China.

919    Mycosphere 7, 1147–1163.

920  Katoh, K., Misawa, K., Kuma, K., Miyata, T., 2002. MAFFT: a novel method for rapid

921    multiple sequence alignment based on fast Fourier transform. Nucleic Acids

922    Research 30, 3059–3066.

923  Katoh, K., Standley, D.M, 2013. MAFFT multiple sequence alignment software

924    version 7: Improvements in performance and usability. Molecular Biology and

925    Evolution 30, 772–780.

926  Kumar, S., Stecher, G., Tamura, K., 2016. MEGA7: Molecular Evolutionary Genetics

927    Analysis version 7.0 for bigger datasets. Molecular Biology and Evolution 33,

928    1870–1874.

929  Larget, B., Kotha, S.K., Dewey, C.N., Ané, C., 2010. BUCKy: Gene tree/species tree

930    reconciliation with the Bayesian concordance analysis. Bioinformatics 26, 2910–

931    2911.

932  Librado, P., Rozas, J.,2009. DnaSP v5: a software for comprehensive analysis of

933    DNA polymorphism data. Bioinformatics 25, 1451–1452.

934  Lima, N.B., Batista, M.V.A., Morais Jr, M.A., Barbosa, M.A.G., Michereff, S.J., Hyde,

935    K.D., Câmara, M.P.S., 2013. Five *Colletotrichum* species are responsible for

936    mango anthracnose in northeastern Brazil. Fungal Diversity 61:75–88.

937  Liu, F., Cai, L., Crous, P.W., Damm, U., 2014. The *Colletotrichum gigasporum*

938    species complex. Persoonia 33, 83–97.

939    Liu, F., Wang, M., Damm, U., Crous, P.W., Cai, L., 2016. Species boundaries in plant

940        pathogenic fungi: A *Colletotrichum* case study. BMC Evolutionary Biology 16, 81.

941    Liu, F., Weir, B.S., Damm, U., Crous, P.W., Wang, Y., Liu, B., Wang, M., Zhang, M.,

942        Cai, L., 2015. Unravelling *Colletotrichum* species associated with *Camellia*:

943        Employing ApMat and GS loci to resolve species in the *C. gloeosporioides*

944        complex. Persoonia 35, 63–86.

945    Lopez-Giraldez, F., Townsend, J.P., 2011. PhyDesign: An online application for

946        profiling phylogenetic informativeness. BMC Evolutionary Biology 11, 152.

947    Magain, N., Miadlikowska, J., Mueller, O., Gajdeczka, M., Truong, C., Salamov, A.A.,

948        Dubchak, I., Grigoriev, I.V., Goffinet, B., Sérusiaux, E., Lutzoni, F., 2017.

949        Conserved genomic collinearity as a source of broadly applicable, fast evolving,

950        markers to resolve species complexes: A case study using the lichen-forming

951        genus *Peltigera* section *Polydactylon*. Molecular Phylogenetics and Evolution

952        117, 10–29.

953    Marin-Felix, Y., Groenewald, J.Z., Cai, L., Chen, Q., Marincowitz, S., Barnes, I.,

954        Bensch, K., Braun, U., Camporesi, E., Damm, U., de Beer, Z.W., Dissanayake,

955        A., Edwards, J., Giraldo, A., Hernandez-Restrepo, M., Hyde, K.D., Jayawardena,

956        R.S., Lombard, L., Luangsa-Ard, J., McTaggart, A.R., Rossman, A.Y., Sandoval-

957        Denis, M., Shen, M., Shivas, R.G., Tan, Y.P., van der Linde, E.J., Wingfield, M.J.,

958        Wood, A.R., Zhang, J.Q., Zhang, Y., Crous, P.W., 2017. Genera of

959        phytopathogenic fungi: GOPHY 1. Studies in Mycology 86, 99–216.

960    Mills, P.R., Hodson, A., Brown, A.E., 1992. Molecular differentiation of *Colletotrichum*

961        *gloeosporioides* isolates infecting tropical fruits, in: Bailey, J.A., Jeger, M.J.

962        (Eds.), *Colletotrichum*: Biology, Pathology and Control, CABI, Wallingford, 269–

963        288.

964    Moriwaki,J., Tsukiboshi,T., 2009. *Colletotrichum echinochloae*, a new species on

965    Japanese barnyard millet (*Echinochloa utilis*). Mycoscience 50, 273–280.

966    Niu, X., Gao, H., Qi, J., Chen, M., Tao, A., Xu, J., Dai, Z., Su, J., 2016. *Colletotrichum*

967    species associated with jute (*Corchorus capsularis* L.) anthracnose in

968    southeastern China. Scientific Reports 6, 25179.

969    O'Connell, R.J., Thon, M.R., Hacquard, S., Amyotte, S.G., Kleemann, J., Torres,

970    M.F., Damm, U., Buiate, E.A., Epstein, L., Alkan, N., Altmuller, J., Alvarado-

971    Balderrama, L., Bauser, C.A., Becker, C., Birren, B.W., Chen, Z., Choi, J.,

972    Crouch, J.A., Duvick, J.P., Farman, M.A., Gan, P., Heiman, D., Henrissat, B.,

973    Howard, R.J., Kabbage, M., Koch, C., Kracher, B., Kubo, Y., Law, A.D., Lebrun,

974    M.H., Lee, Y.H., Miyara, I., Moore, N., Neumann, U., Nordstrom, K., Panaccione,

975    D.G., Panstruga, R., Place, M., Proctor, R.H., Prusky, D., Rech, G., Reinhardt,

976    R., Rollins, J.A., Rounsley, S., Schardl, C.L., Schwartz, D.C., Shenoy, N.,

977    Shirasu, K., Sikhakolli, U.R., Stuber, K., Sukno, S.A., Sweigard, J.A., Takano, Y.,

978    Takahara, H., Trail, F., Van Der Does, H.C., Voll, L.M., Will, I., Young, S., Zeng,

979    Q., Zhang, J., Zhou, S., Dickman, M.B., Schulze-Lefert, P., Van Themaat, E.V.,

980    Ma, L.J., Vaillancourt,L.J., 2012. Lifestyle transitions in plant pathogenic

981    *Colletotrichum* fungi deciphered by genome and transcriptome analyses. Nature

982    Genetics 44, 1060–1065.

983    Oliveira, L.F.M., Feijó, F.M., Mendes, A.L.S.F., Neto, J.D.V., Netto, M.S.B.,

984    Assunção, I.P., Lima, G.S.A., 2018. Identification of *Colletotrichum* species

985    associated with brown spot of cactus prickly pear in Brazil. Tropicala Plant

986    Pathology 43, 247–253.

987    Paradis, E., Claude, J., Strimmer, K., 2004. APE: Analyses of phylogenetics and

988    evolution in R language. Bioinformatics 20, 289–290.

989    Pond, S.L., Frost, S.D., Muse, S.V., 2005. HyPhy: hypothesis testing using

990        phylogenies. Bioinformatics 21, 676–679.

991    R Core Team. 2017. R: A language and environment for statistical computing. R

992        Foundation for Statistical Computing, Vienna, Austria.

993    Rojas, E.I., Rehner, S.A., Samuels, G.J., Van Bael, S.A., Herre, E.A., Cannon, P.,

994        Chen, R., Pang, J., Wang, R., Zhang, Y., Peng, Y.Q., Sha, T., 2010.

995        *Colletotrichum gloeosporioides s.l.* associated with *Theobroma cacao* and other

996        plants in Panamá: Multilocus phylogenies distinguish host-associated pathogens

997        from asymptomatic endophytes. Mycologia 102,1318–1338.

998    Ronquist, F., Teslenko, M., van der Mark, P., Ayres, D.L., Darling, A., Höhna, S.,

999        Larget, B., Liu, L., Suchard, M.A., Huelsenbeck, J.P., 2012. MrBayes v. 3.2:

1000       Efficient Bayesian phylogenetic inference and model choice across a large model

1001       space. Systematic Biology 61, 539–542.

1002   Sakalidis, M.L., Hardy, G.E.S.J., Burgess, T.I., 2011. Use of the Genealogical Sorting

1003       Index (GSI) to delineate species boundaries in the *Neofusicoccum parvum-*

1004       *Neofusicoccum ribis* species complex. Molecular Phylogenetics and Evolution,

1005       60, 333–344.

1006   Samarakoon, M.C., Peršoh, D., Hyde, K.D., Bulgakov, T.S., Manawasinghe, I.S.,

1007       Jayawardena, R.S., Promputtha, I., 2018. *Colletotrichum acidae* sp. nov. from

1008       northern Thailand and a new record of *C. dematium* on Iris sp. Mycosphere 9:

1009       583–597.

1010   Schmitt,.I, Crespo, A., Divakar, P.K., Fankhauser, J.D., Herman-Sackett, E., Kalb, K.,

1011       Nelsen, M.P., Nelson, N.A., Rivas-Plata, E., Shimp, A.D.,Widhelm, T., Lumbsch,

1012       H.T., 2009. New primers for promising single-copy genes in fungal phylogenetics

1013       and systematics. Persoonia 23, 35–40.

1014    Sela, I., Ashkenazy, H., Katoh, K., Pupko, T., 2015. GUIDANCE2: Accurate detection

1015        of unreliable alignment regions accounting for the uncertainty of multiple

1016        parameters. Nucleic Acids Research 43, W7–W14.

1017    Sharma, G., Kumar, N., Weir, B.S., Hyde, K.D., Shenoy, B.D., 2013. The ApMat

1018        marker can resolve *Colletotrichum* species: A case study with *Mangifera indica*.

1019        Fungal Diversity 61,117–38.

1020    Sharma, G., Pinnaka, A.K., Belle, D.S., 2015. Resolving the *Colletotrichum siamense*

1021        species complex using ApMat marker. Fungal Diversity 71, 247–64.

1022    Sharma, G., Maymon, M., Freeman, S., 2017. Epidemiology, pathology and

1023        identification of *Colletotrichum* including a novel species associated with avocado

1024        (*Persea americana*) anthracnose in Israel. Scientific Reports 7, 15839.

1025    Silva, D.N., Talhinas, P., Várzea, V., Cai, L., Paulo, O.S., Batista, D., 2012.

1026        Application of the Apn2/MAT locus to improve the systematics of the

1027        *Colletotrichum gloeosporioides* complex: An example from coffee (*Coffea* spp.)

1028        hosts. Mycologia 104, 396–409.

1029    Silva, J.R.A, Chaves, T.P., Silva,a A.R.G., Barbosa, L.F., Costa, F.F.O., Ramos-

1030        Sobrinho, R., Teixeira, R.R.O., Silva, S.J.C., Lima, G.S.A., Assunção, I.P., 2017.

1031        Molecular and morpho-cultural characterization of *Colletotrichum* spp. associated

1032        with anthracnose on *Capsicum* spp. in northeastern Brazil. Tropical Plant

1033        Pathology 42, 315–319.

1034    Sousa, E.S., Silva, J.R.A., Assunção, I.P., Melo, M.P., Feijó, F.M., Matos, K.S., Lima,

1035        G.S.A., Beserra Jr, J.E.A., 2018. *Colletotrichum* species causing anthracnose on

1036        lima bean in Brazil. Tropical Plant Pathology 43, 78–84.

1037   Sreenivasaprasad, S., Brown, A.E., Mills, P.R., 1992. DNA sequence variation and

1038       interrelationship among *Colletotrichum* species causing strawberry anthracnose.

1039       Physiological and Molecular Plant Pathology 41, 265–281.

1040   Stawatakis, A., 2014. RAxML version 8: A tool for Phylogenetic Analysis and Post-

1041       Analysis of Large Phylogenies. Bioinformatics 30, 1312–1313.

1042   Sutton, B.C., 1980. The Coelomycetes. Fungi Imperfecti with Pycnidia, Acervuli and

1043       Stromata. CABI, Kew

1044   Tao, G., Liu, Z.Y., Liu, F., Gao, Y.H., Cai,L., 2013. Endophytic *Colletotrichum* species

1045       from *Bletilla ochracea* (Orchidaceae), with descriptions of seven new speices.

1046       Fungal Diversity 61, 139–164.

1047   Taylor, J.W., Jacobson, D.J., Kroken, S., Kasuga, T., Geiser, D.M., Hibbett, D.S.,

1048       Fisher, M.C., 2000. Phylogenetic species recognition and species concepts in

1049       fungi. Fungal Genetics and Biology 31, 21–32.

1050   Vaidya, G., Lohman, D. J., Meier, R., 2011. SequenceMatrix: concatenation software

1051       for the fast assembly of multi-gene datasets with character set and codon

1052       information. Cladistics 27, 171–180.

1053   Veloso, J.S., Câmara, M.P.S., Lima, W.G., Michereff, S.J., Doyle, V.P.. (2018). Why

1054       species delimitation matters for fungal ecology: *Colletotrichum* diversity on wild

1055       and cultivated cashew in Brazil. Fungal Biology 122, 677–691.

1056   Vieira, W.A., Michereff, S.J., de Morais Jr, M.A., Hyde, K.D., Câmara, M.P.S., 2014.

1057       Endophytic species of *Colletotrichum* associated with mango in northeastern

1058       Brazil. Fungal Diversity 67, 181–202.

1059   Vieira, W.A.S., Lima, W.G., Nascimento, E.S., Michereff, S.J., Câmara, M.P.S.,

1060       Doyle, V.P., 2017. The impact of phenotypic and molecular data on the inference

1061       of *Colletotrichum* diversity associated with *Musa*. Mycologia 109, 912–934.

1062    von Arx, J.A., 1957. Die Arten der Gattung *Colletotrichum* Cda. Phytopathologische

1063        Zeitschrift 29, 413–468.

1064    Wang, Q.T., Liu, X.T., Ma, H.Y., Shen, X.Y., Hou, C.L., 2019. *Colletotrichum*

1065        *yulongense sp. nov.* and *C. rhombiforme* isolated as endophytes from *Vaccinium*

1066        *dunalianum var. urophyllum* in China. Phytotaxa 394, 285–298.

1067    Weir, B.S., Johnston, P.R., Damm, U., 2012. The *Colletotrichum gloeosporioides*

1068        species complex. Studies in Mycology 73,115–180.

**Table 1**. Alignment and phylogenetic informativeness profile statistics for markers used in *Colletotrichum* by species complex.

**Table 1**.

| C. acutatum s. l. | Length[1] | Invariable[1] | Variable[1] | Singletons[1] | Parsymony informative[1] | PIV ($10^{-6}$)[2] | Pimax[3] |
|---|---|---|---|---|---|---|---|
| ACT | 239 | 163 | 76 | 24 | 52 | 15.36 | 0.99 |
| CHS-1 | 251 | 216 | 35 | 5 | 30 | 8.56 | 0.78 |
| GAPDH | 272 | 165 | 94 | 18 | 76 | 19.95 | 0.78 |
| HIS3 | 386 | 298 | 88 | 18 | 70 | 24.96 | 0.67 |
| nrITS | 454 | 427 | 26 | 7 | 19 | 5.41 | 0.99 |
| TUB2 | 481 | 386 | 95 | 24 | 71 | 19.7 | 0.99 |
| *C. boninense s. l.* | | | | | | | |
| ACT | 277 | 199 | 78 | 15 | 63 | 29.07 | 0.49 |
| CAL | 438 | 293 | 142 | 16 | 126 | 52.86 | 0.58 |
| CHS-1 | 277 | 228 | 49 | 4 | 45 | 22.39 | 0.48 |
| GAPDH | 247 | 163 | 83 | 9 | 74 | 47.6 | 0.24 |
| HIS3 | 389 | 292 | 93 | 11 | 82 | 46.96 | 0.32 |
| nrITS | 537 | 490 | 47 | 11 | 36 | 18.41 | 0.39 |
| TUB2 | 483 | 345 | 137 | 23 | 114 | 48.35 | 0.63 |
| *C. dematium s. l.* | | | | | | | |
| ACT | 235 | 188 | 46 | 1 | 45 | 44.65 | 0.09 |
| CHS-1 | 251 | 215 | 36 | 6 | 30 | 21.48 | 0.5 |
| GAPDH | 264 | 140 | 122 | 13 | 109 | 237.56 | 0.05 |
| HIS3 | 371 | 301 | 70 | 7 | 63 | 183.1 | 0.05 |
| nrITS | 517 | 492 | 25 | 1 | 24 | 32.94 | 0.06 |
| TUB2 | 497 | 389 | 105 | 10 | 95 | 118.42 | 0.06 |
| *C. destructivum s. l.* | | | | | | | |
| ACT | 263 | 218 | 44 | 20 | 24 | 18 | 0.66 |
| CHS-1 | 280 | 257 | 23 | 11 | 12 | 10.13 | 0.63 |
| GAPDH | 225 | 149 | 47 | 16 | 31 | 27.35 | 0.28 |
| HIS3 | 389 | 318 | 71 | 18 | 53 | 38.75 | 0.37 |
| nrITS | 557 | 526 | 26 | 13 | 13 | 12.41 | 0.17 |
| TUB2 | 514 | 432 | 77 | 33 | 44 | 32.49 | 0.73 |
| *C. dracaenophilum s. l.* | | | | | | | |
| ACT | 254 | 205 | 49 | 6 | 43 | 31.72 | 0.44 |
| CHS-1 | 282 | 257 | 25 | 5 | 20 | 18.53 | 0.46 |
| GAPDH | 273 | 192 | 81 | 14 | 67 | 52.08 | 0.41 |
| HIS3 | 417 | 354 | 62 | 7 | 55 | 35.74 | 0.45 |
| nrITS | 541 | 504 | 37 | 13 | 24 | 28.58 | 0.39 |
| TUB2 | 492 | 411 | 81 | 9 | 72 | 58.22 | 0.4 |
| *C. gigasporum s. l.* | | | | | | | |
| ACT | 280 | 239 | 37 | 17 | 20 | 23.82 | 0.19 |
| CAL | 675 | 543 | 131 | 56 | 75 | 63.14 | 0.44 |
| CHS-1 | 299 | 266 | 33 | 15 | 18 | 26.38 | 0.11 |
| GAPDH | 290 | 193 | 94 | 51 | 43 | 61.02 | 0.28 |
| GS | 702 | 521 | 176 | 92 | 84 | 105.17 | 0.10 |
| HIS3 | 416 | 362 | 53 | 22 | 31 | 45.36 | 0.12 |
| nrITS | 545 | 500 | 39 | 15 | 24 | 34.85 | 0.11 |
| TUB2 | 537 | 451 | 86 | 40 | 46 | 58.50 | 0.15 |

**C. gloeosporioides s. l.**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| APN2 | 735 | 624 | 108 | 18 | 90 | 73.91 | 0.17 |
| APN2/MAT-IGS | 579 | 378 | 201 | 4 | 197 | 149.81 | 0.19 |
| CAL | 648 | 578 | 70 | 1 | 69 | 48.76 | 0.27 |
| GAP2-IGS | 708 | 569 | 135 | 14 | 121 | 79.71 | 0.42 |
| GAPDH | 767 | 734 | 33 | 2 | 31 | 24.51 | 0.33 |
| GS | 665 | 571 | 87 | 24 | 63 | 43.83 | 0.39 |
| TUB2 | 1231 | 1116 | 115 | 10 | 105 | 67.23 | 0.49 |

**C. graminicola s. l.**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ACT | 269 | 174 | 82 | 52 | 30 | 30.25 | 0.68 |
| CHS-1 | 280 | 237 | 43 | 24 | 19 | 17.06 | 0.55 |
| nrITS | 454 | 375 | 70 | 37 | 33 | 30.38 | 0.23 |
| TUB2 | 500 | 350 | 145 | 89 | 56 | 58.96 | 0.46 |

**C. magnum s. l.**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ACT | 276 | 260 | 16 | 13 | 3 | 11.16 | 0.40 |
| CHS-1 | 257 | 247 | 10 | 7 | 3 | 11.16 | 0.40 |
| GAPDH | 241 | 215 | 26 | 14 | 12 | 13.35 | 0.42 |
| HIS3 | 399 | 360 | 39 | 28 | 11 | 26 | 0.44 |
| nrITS | 539 | 533 | 6 | 6 | 0 | 4.16 | 0.40 |
| TUB2 | 536 | 500 | 36 | 24 | 12 | 23.15 | 0.52 |

**C. orbiculare s. l.**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ACT | 226 | 202 | 21 | 5 | 16 | 18.37 | 0.36 |
| CHS-1 | 280 | 268 | 12 | 2 | 10 | 6.16 | 0.65 |
| GAPDH | 243 | 203 | 40 | 8 | 32 | 32.01 | 0.17 |
| GS | 954 | 810 | 134 | 41 | 93 | 69.13 | 0.53 |
| HIS3 | 386 | 334 | 50 | 10 | 40 | 36.51 | 0.13 |
| nrITS | 529 | 514 | 15 | 4 | 11 | 6.38 | 0.52 |
| TUB2 | 467 | 431 | 36 | 9 | 27 | 19.75 | 0.55 |

**C. orchidearum s. l.**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ACT | 276 | 244 | 29 | 17 | 12 | 7.06 | 0.52 |
| CHS-1 | 265 | 243 | 22 | 6 | 16 | 11.46 | 0.47 |
| GAPDH | 242 | 195 | 47 | 37 | 10 | 7.99 | 0.77 |
| HIS3 | 402 | 343 | 59 | 31 | 28 | 16.48 | 0.67 |
| nrITS | 539 | 521 | 18 | 10 | 8 | 5.81 | 0.57 |
| TUB2 | 540 | 460 | 79 | 51 | 28 | 17.44 | 0.80 |

**C. spaethianum s. l.**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ACT | 210 | 152 | 56 | 28 | 28 | 19.44 | 0.77 |
| GAPDH | 240 | 132 | 105 | 28 | 77 | 32.53 | 0.47 |
| nrITS | 482 | 470 | 10 | 4 | 6 | 8.66 | 0.73 |
| TUB2 | 443 | 351 | 92 | 44 | 48 | 33.26 | 0.62 |

**C. truncatum s. l.**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ACT | 215 | 179 | 36 | 6 | 30 | 22.89 | 0.41 |
| GAPDH | 276 | 186 | 90 | 6 | 84 | 56.86 | 0.46 |
| nrITS | 479 | 459 | 20 | 4 | 16 | 10.17 | 0.53 |
| TUB2 | 434 | 353 | 76 | 9 | 67 | 33.14 | 0.53 |

1 - number of base pairs.
2 - PIV, phylogenetic informativeness values on a per net basis
3 - PImax, optimal divergence time. Values are expressed in arbitrary units.

1072 **FIGURE CAPTIONS**

1073 **Fig. 1.** Barcode gap and distance overlap between the intra- and inter-specific

1074 distances. Values were calculated based on the intra- and inter-specific distances

1075 frequencies distribution of each *Colletotrichum* species complex.

1076

1077 **Fig. 2.** Ultrametric trees and net phylogenetic informativeness profiles of markers

1078 used for phylogenetic studies of 13 *Colletotrichum* species complexes. Values on the

1079 X-axes correspond to the relative timescale (0—1) based on the root-to-tip distance.

1080 Values on the Y-axes represent net phylogenetic informativeness values ($10^{-6}$) in

1081 arbitrary units.

1082

1083 **Fig. 3.** Heat map of the Genealogical Sorting Indices (GSI) by *Colletotrichum* species

1084 complex. GSIs of 1000 bootstrap trees were calculated with 100 permutations. Rows

1085 correspond to species, and columns correspond to individual markers and

1086 concatenated datasets (all markers and best markers). Asterisks represent the best

1087 markers for each complex.

1088

1089 **Fig. 4.** Primary concordance trees resulting from the Bayesian concordance analyses

1090 including isolates from the *C. acutatum* complex. A. All markers (ACT, CHS-1,

1091 GAPDH, HIS3, ITS and TUB2). B. Best markers (GAPDH, HIS3 and TUB2).

1092 Concordance factors are shown above the branches that were resolved by at least

1093 one marker (≥0.16 for all markers and ≥0.33 for the best markers).

1094

1095 **Fig. 5.** Primary concordance trees resulting from the Bayesian concordance analyses

1096 including isolates from *C. boninense* complex. A. All markers (ACT, CAL, CHS-1,

1097 GAPDH, HIS3, ITS and TUB2). B. Best markers (CAL, GAPDH and TUB2).

1098 Concordance factors are shown above the branches that were resolved by at least

1099 one marker (≥0.14 for all markers and ≥0.33 for the best markers).

1100

1101 **Fig. 6.** Primary concordance trees resulting from the Bayesian concordance analyses

1102 including isolates from the *C. dematium* complex. A. All markers (ACT, CHS-1,

1103 GAPDH, HIS3, ITS and TUB2). B. Best markers (GAPDH, HIS3 and TUB2).

1104 Concordance factors are shown above the branches that were resolved by at least

1105 one marker (≥0.16 for all markers and ≥0.33 for the best markers).

1106

1107 **Fig. 7.** Primary concordance trees resulting from the Bayesian concordance analyses

1108 including isolates from the *C. destructivum* complex. A. All markers (ACT, CHS-1,

1109 GAPDH, HIS3, ITS and TUB2). B. Best markers (GAPDH, HIS3 and TUB2).

1110 Concordance factors are shown above the branches that were resolved by at least

1111 one marker (≥0.16 for all markers and ≥0.33 for the best markers).

1112

1113 **Fig. 8.** Primary concordance trees resulting from the Bayesian concordance analyses

1114 including isolates from the *C. dracaenophilum* complex. A. All markers (ACT, CHS-1,

1115 GAPDH, HIS3, ITS and TUB2). B. Best markers (GAPDH, HIS3 and TUB2).

1116 Concordance factors are shown above the branches that were resolved by at least

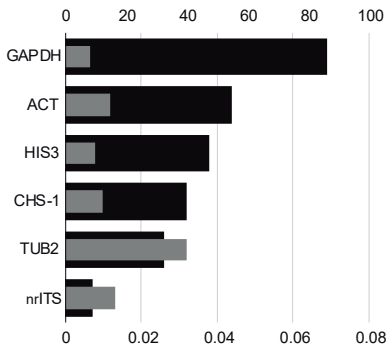1117 one marker (≥0.16 for all markers and ≥0.33 for the best markers).

1118

1119 **Fig. 9.** Primary concordance trees resulting from the Bayesian concordance analyses

1120 including isolates from the *C. gigasporum* complex. A. All markers (ACT, CAL, CHS-

1121 1, GAPDH, GS, HIS3, ITS and TUB2). B. Best markers (CAL, GAPDH and GS).

1122   Concordance factors are shown above the branches that were resolved by at least

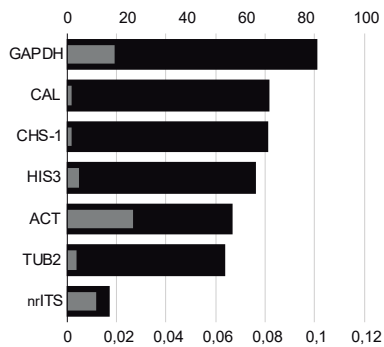1123   one marker (≥0.13 for all markers and ≥0.33 for the best markers).

1124

1125   **Fig. 10.** Primary concordance trees resulting from the Bayesian concordance

1126   analyses including isolates from the *C. gloeosporioides* complex. A. All markers

1127   (APN2, APN2/MAT-IGS, CAL, GAPDH, GAP2-IGS, GS and TUB2). B. Best markers

1128   (APN2, APN2/MAT-IGS and GAP2-IGS). Concordance factors are shown above the

1129   branches that were resolved by at least one marker (≥0.14 for all markers and ≥0.33

1130   for the best markers).

1131

1132   **Fig. 11.** Primary concordance trees resulting from the Bayesian concordance

1133   analyses including isolates from the *C. magnum* complex. A. All markers (ACT, CHS-

1134   1, GAPDH, HIS3, ITS and TUB2). B. Best markers (GAPDH, HIS3 and TUB2).

1135   Concordance factors are shown above the branches that were resolved by at least

1136   one marker (≥0.16 for all markers and ≥0.33 for the best markers).

1137

1138   **Fig. 12.** Primary concordance trees resulting from the Bayesian concordance

1139   analyses including isolates from the *C. orbiculare* complex. A. All markers (ACT,

1140   CHS-1, GAPDH, GS, HIS3, ITS and TUB2). B. Best markers (HIS3, GAPDH and

1141   GS). Concordance factors are shown above the branches that were resolved by at

1142   least one marker (≥0.14 for all markers and ≥0.33 for the best markers).

1143

1144   **Fig. 13.** Primary concordance trees resulting from the Bayesian concordance

1145   analyses including isolates from the *C. orchidearum* complex. A. All markers (ACT,

1146   CHS-1, GAPDH, HIS3, ITS and TUB2). B. Best markers (CHS-1, HIS3 and TUB2).

1147    Concordance factors are shown above the branches that were resolved by at least

1148    one marker (≥0.16 for all markers and ≥0.33 for the best markers).

1149

1150    **Fig. 14.** Primary concordance trees resulting from the Bayesian concordance

1151    analyses including isolates from the *C. spaethianum* complex. A. All markers (ACT,

1152    GAPDH, ITS and TUB2). B. Best markers (ACT, GAPDH, and TUB2). Concordance

1153    factors are shown above the branches that were resolved by at least one marker

1154    (≥0.25 for all markers and ≥0.33 for the best markers).

1155

1156    **Fig. 15.** Primary concordance trees resulting from the Bayesian concordance

1157    analyses including isolates from the *C. truncatum* complex. A. All markers (ACT,

1158    GAPDH, ITS and TUB2). B. Best markers (ACT, GAPDH, and TUB2). Concordance

1159    factors are shown above the branches that were resolved by at least one marker

1160    (≥0.25 for all markers and ≥0.33 for the best markers).

1161

1162

1163

Figure showing barcode gap analysis for multiple *Colletotrichum* species complexes (*C. acutatum s. l.*, *C. boninense s. l.*, *C. dematium s. l.*, *C. destructivum s. l.*, *C. dracaenophilum s. l.*,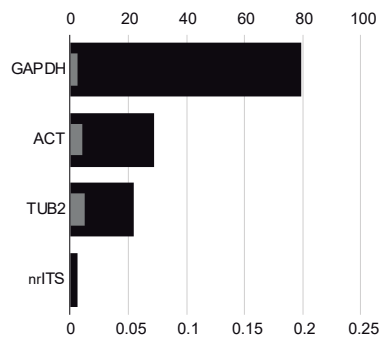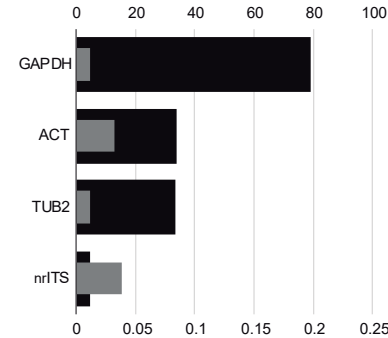 *C. gigasporum s. l.*, *C. gloeosporioides s. l.*, *C. magnum s. l.*, *C. orbiculare s. l.*, *C. orchidearum s. l.*, *C. spaethianum s. l.*, *C. truncatum s. l.*).
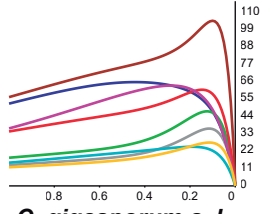
**Barcode gap** (black) — **Intra/inter-specific distance overlap** (gray)

*C. acutatum s. l.*

*C. boninense s. l.*

*C. destructivum s. l.*

*C. graminicola s. l.*

*C. spaethianum s. l.*

*C. orbiculare s. l.*

*C. orchidearum s. l.*

*C. gloeosporioides s. l.*

*C. dematium s. l.*

*C. magnum s. l.*

*C. gigasporum s. l.*

*C. truncatum s. l.*

*C. dracaenophilum s. l.*

ACT    CAL    GAP2-IGS    nrITS

APN2    CHS-1    GS    TUB2

APN2/MAT-IGS    GAPDH    HIS3

**C. truncatum s. l.**

**C. dracaenophilum s. l.**

**C. gigasporum s. l.**

**C. magnum s. l.**

**C. orchidearum s. l.**

**C. orbiculare s. l.**

**C. dematium s. l.**

**C. spaethianum s. l.**

**C. acutatum s. l.**

**C. boninense s. l.**

**C. destructivum s. l.**

**C. gloeosporioides s. l.**

GSI values

0    0.25    0.50    0.75    1

A

B

Clade 2
0.37

| | 0.48 |
| 0.37 |
| *C. scovillei* |
| *C. guajavae* |

(Phylogenetic tree figure)

Panel A (left tree):
- 0.48 / 0.37 *C. scovillei*
- *C. guajavae*
- 0.17 *C. cairnsense*
- 0.58 *C. brisbanense*
- 0.21 / 0.39 / 0.19 *C. laticiphilum*
- *C. indonesiense*
- 0.44 / 0.18 *C. cosmi*
- *C. walleri*
- 0.66 *C. nymphaeae*
- 0.42 / 0.38 *C. paxtonii*
- 0.41 / 0.17 *C. simmondsii*
- *C. sloanei*
- 0.84 *C. chrysanthemi*
- 0.21 (Clade 2, 0.37)
- 0.31 *C. paranaense*
- 0.37 / 0.18 *C. limeticola*
- *C. melonis*
- 0.38 *C. costaricense*
- 0.17 *C. lupini*
- 0.74 *C. tamarilloi*
- 0.18 / 1 *C. cuscutae*
- Clade 1, 0.83
- 0.59 *C. abscissum*
- Clade 3, 1
- *C. fioriniae*
- 1 *C. acerbum*
- 0.42 / 0.23 *C. rhombiforme*
- 0.28 *C. phormii*
- 0.43 *C. kinghornii*
- 0.84 / 0.27 *C. australe*
- *C. salicis*
- 0.95 *C. johnstonii*
- 0.69 / 0.51 *C. pyricola*
- *C. godetiae*
- 0.85 / 0.26 / 0.18 (Clade 5)
- Clade 4, 0.18 *C. acutatum*

Panel B (right tree):
- 0.73 / 0.4 *C. scovillei*
- *C. guajavae*
- 1 *C. chrysanthemi*
- *C. cairnsense*
- 0.33 / 0.55 / 0.36 *C. cosmi*
- *C. walleri*
- 0.99 *C. nymphaeae*
- 0.73 / 0.68 / 0.73 / 0.34 *C. paxtonii*
- *C. simmondsii*
- *C. sloanei*
- 0.33 / 0.72 / 0.37 *C. laticiphilum*
- *C. indonesiense*
- *C. brisbanense*
- Clade 2, 0.74
- 0.73 / 0.33 / 0.63 / 0.37 *C. abscissum*
- *C. costaricense*
- 1 / 0.35 *C. tamarilloi*
- 0.78 *C. lupini*
- *C. paranaense*
- *C. limeticola*
- 0.36 / 0.38 *C. melonis*
- *C. cuscutae*
- Clade 1, 1
- Clade 3, 1
- *C. fioriniae*
- 0.33 / 0.33
- 1 *C. acerbum*
- 0.39 / 0.44 *C. rhombiforme*
- 0.5 *C. phormii*
- 0.49 *C. kinghornii*
- 1 / 0.41 *C. australe*
- *C. salicis*
- 0.9 *C. johnstonii*
- 0.98 / 0.67 *C. pyricola*
- *C. godetiae*
- Clade 5
- 0.42
- 1 *C. acutatum* — Clade 4

A

| 0.45 | | *C. cymbidiicola* |
| 0.73 | 0.34 | *C. oncidii* |
| 0.68 | 0.85 | *C. boninense* |
| 0.4 | 0.78 | *C. torulosum* |
| 0.85 | 0.79 | *C. colombiense* |
| 0.38 | 0.79 | *C. beeveri* |
| | | *C. brassicicola* |
| 0.88 | | *C. brasiliense* |
| 0.45 | 0.99 | *C. parsoniae* |
| 0.4 | 1 | *C. hippeastri* |
| 0.91 | 1 | *C. constrictum* |
| 0.21 | | *C. dacrycarpi* |
| 0.57 | | *C. catinaense* |
| 0.49 | 0.19 | *C. limonicola* |
| 0.99 | 0.26 | *C. novae-zelandiae* |
| 0.16 | 0.2 | *C. karstii* |
| 0.16 | 0.72 | |
| 0.49 | 0.74 | *C. phyllanti* |
| | | *C. annelatum* |
| 0.81 | 1 | *C. petchii* |

B

| *C. cymbidiicola* | 0.53 |
| *C. oncidii* | 0.46 | 1 |
| *C. boninense* | 1 | 0.76 |
| *C. torulosum* | 0.56 | 0.77 |
| *C. colombiense* | 1 | 1 |
| *C. beeveri* | 0.46 | 0.77 |
| *C. brassicicola* | 0.37 |
| *C. constrictum* | 0.78 |
| *C. dacrycarpi* | 1 |
| *C. brasiliense* | 1 |
| *C. hippeastri* | 0.52 | 1 | 1 |
| *C. parsoniae* | |
| *C. karstii* | |
| *C. phyllanti* | 0.69 | 0.78 |
| *C. annelatum* | 1 | 1 |
| *C. petchii* | 1 |
| *C. catinaense* | 0.5 |
| *C. limonicola* | 0.36 | 0.45 |
| *C. novae-zelandiae* | 0.98 | 0.38 |

A

- 0.49
- 0.5
- 0.56
- 1
- 0.53
- *C. dematium*
- 0.35
- 0.49
- *C. lineola*
- 1
- *C. anthrisci*
- *C. fructi*
- 0.56
- 0.54
- 1
- *C. circinans*
- 0.58
- 0.39
- 0.46
- *C. spinaceae*
- 0.49

B

- *C. dematium*
- 0.62
- 0.69
- *C. lineola*
- 0.61
- 1
- *C. anthrisci*
- 1
- 0.58
- *C. circinans*
- 0.73
- 0.35
- 1
- *C. spinaceae*
- 1
- 1
- *C. fructi*

A

| Label | Species |
|---|---|
| 0.35 | *C. higginsianum* |
| 0.61 | *C. vignae* |
| 0.18 | *C. fuscum* |
| 0.29 / 0.23 / 0.18 | *C. antirrhinicola* |
| 0.36 | *C. bryoniicola* |
| 0.37 | *C. utrechtense* |
| 0.27 / 0.71 | *C. tabacum* |
| 0.34 / 0.29 / 0.99 | *C. lentis* |
| 0.37 | *C. pisicola* |
| 1 | *C. lini* |
| 0.19 / 0.28 / 0.53 / 0.28 | *C. americae* |
| 0.42 / 0.83 | *C. destructivum* |
| 0.29 / 0.99 / 0.99 | *C. ocimi* |
| 0.2 | |

B

| Species | Label |
|---|---|
| *C. fuscum* | 0.4 |
| *C. vignae* | 0.34 / 0.69 |
| *C. bryoniicola* | |
| *C. antirrhinicola* | 0.67 |
| *C. higginsianum* | 0.35 |
| *C. utrechtense* | 0.41 / 0.67 |
| *C. tabacum* | 0.49 / 0.39 / 0.98 |
| *C. lini* | 0.41 / 0.39 |
| *C. americae* | 1 |
| *C. destructivum* | 0.55 / 0.77 / 1 |
| *C. ocimi* | |
| *C. lentis* | 1 / 0.37 |
| *C. pisicola* | |
| | 0.4 |

A

0.41

1

1

0.88

1

1

B

*C. dracaenophilum*

*C. yunnanense*

*C. coelogynes*

*C. tropicicola*

1

1

1

1

A

- 0.38
  - 1
    - 0.5
      - *C. gigasporum*
    - *C. pseudomajus*
  - 0.51
    - 0.37
      - 1
        - *C. vietnamense*
        - *C. arxii*
      - *C. radicis*

B

- 0.63
  - 0.73
    - 1
      - *C. gigasporum*
    - *C. pseudomajus*
  - 0.37
    - 1
      - 1
        - *C. vietnamense*
        - *C. radicis*
      - *C. arxii*

A

- 0.52
- 1 *C. theobromicola*
- 0.4
- 0.4
- 1 *C. musae*
- 0.27
- 0.26
- 0.3 *C. siamense*
- 0.59
- 0.3
- 0.71
- 0.85 *C. queenslandicum*
- 0.27
- 0.99 *C. tropicale*
- 0.21
- 0.3 *C. chrysophilum*
- 0.88
- 0.9
- 0.9 *C. fructicola*

B

- 0.67
- 0.38 *C. siamense*
- 0.65
- 0.33
- 0.77
- 1 *C. queenslandicum*
- 0.56
- 1 *C.theobromicola*
- 1 *C. tropicale*
- 0.71
- 1 *C. chrysophilum*
- 0.78
- *C. fructicola*
- 0.45
- 1 *C. musae*

A

- 0.43
  - 0.55
    - 0.23
      - 0.33
        - *C. brevisporum*
    - *C. lobatum*
    - *C. merremiae*
  - 0.47
    - 0.72
      - 0.28
        - *C. magnum*
      - *C. panamense*
    - *C. cacao*

B

- 0.51
  - 0.67
    - 0.4
      - 0.46
        - *C. brevisporum*
      - *C. lobatum*
    - *C. merremiae*
  - 0.5
    - 1
      - 0.39
        - *C. magnum*
      - *C. panamense*
    - *C. cacao*

A

- 0.52
- 0.33 — C. orbiculare
- 0.27
- 0.44
- 0.17
- 0.24 — C. sidae
- 0.36 — C. tebeestii
- 0.65 — C. spinosum
- 0.27
- 0.72 — C. malvarum
- 0.88
- 0.21
- 0.39 — C. trifolii
- 0.18
- 0.17
- 0.62
- 0.32 — C. lindemuthianum
- 0.32
- 0.44
- 1 — C. bidentis

B

- C. sidae
- 0.66
- C. tebeestii
- 0.33
- 0.79 — C. spinosum
- 0.38
- 0.39
- 0.44 — C. orbiculare
- 0.59
- 0.34
- 1 — C. malvarum
- 0.73
- 0.47 — C. trifolii
- 0.62
- 0.71
- 1 — C. lindemuthianum
- 0.33 — C. bidentis

A

B

| | |
|---|---|
| 0.33 | |
| 0.71 | |
| 0.5 | |

*C. cattleyicola*
*C. piperis*
*C. sojae*
*C. musicola*
*C. orchidearum*
*C. vittalense*
*C. plurivorum*
*C. cliviicola*

A                                                                                    B

0.55
0.54
C. spaethianum
0.37
0.99
C. lilii
0.75
0.5
0.27
C. guizhouensis
0.75
0.56
C. incanum
0.94
C. riograndense
0.6
0.32
C. tofieldiae
0.99
0.82
C. verruculosum
0.99
0.58
0.28
C. liriopes
0.62
C. bletilum

0.71
0.7
0.5
0.99
1
0.66
0.36
0.64
1
0.54
0.73
0.35
0.98
0.77
0.72
0.35
1
0.77
0.46

A

- 0.55
- 0.25
- *C. corchorum-capsularis*
- 0.35
- 0.25
- *C. truncatum*
- 0.29
- 0.97
- 1  *C. curcumae*
- 0.83  *C. acidae*
- 0.98  *C. fusiforme*

B

- 0.39
- 0.33
- *C. truncatum*
- 0.33
- 0.4
- *C. corchorum-capsularis*
- 0.36
- *C. truncatum*
- 0.96
- 1  *C. curcumae*
- 1  *C. acidae*
- 0.9  *C. fusiforme*