1    Manuscript type: Article

2

**3    Speciation in *Howea* palms occurred in sympatry, was preceded by ancestral**

**4    admixture, and was associated with edaphic and phenological adaptation**

5

6    Owen G. Osborne[a,b], Adam Ciezarek[a], Trevor Wilson[c], Darren Crayn[d], Ian Hutton[e], William J.

7    Baker[f], Colin G.N. Turnbull[g], Vincent Savolainen[a,f,1]

8

9    [a]Department of Life Sciences, Silwood Park Campus, Imperial College London, Ascot SL5

10   7PY, UK; [b]current address: Molecular Ecology and Fisheries Genetics Laboratory, School of

11   Natural Sciences, Bangor University, Bangor UK; [c]Royal Botanic Gardens Sydney, Mrs

12   Macquaries Road, Sydney NSW 2000, Australia; [d]James Cook University, Australian

13   Tropical Herbarium, Sir Robert Norman Building, McGregor Road, Smithfield, QLD 4878,

14   Australia; [e]Lord Howe Island Museum, NSW 2898, Lord Howe Island, Australia; [f]Royal

15   Botanic Gardens Kew, Richmond, Surrey TW9 3AB, UK; [g]Department of Life Sciences,

16   South Kensington Campus, Imperial College London, London SW7 2AZ, UK

17

18

19   [1]Corresponding Author: Vincent Savolainen; Tel +442075942374; Email

20   v.savolainen@imperial.ac.uk

21

**Abstract**

22    **Abstract**

23

24    *Howea* palms are viewed as one of the most clear-cut cases of speciation in sympatry. The

25    sister species *H. belmoreana* and *H. forsteriana* are endemic to the oceanic Lord Howe

26    Island, Australia, where they have overlapping distributions and are reproductively isolated

27    mainly by flowering time differences. However, the potential role of introgression from

28    Australian mainland relatives had not previously been investigated, a process that has

29    recently put other examples of sympatric speciation into question. Furthermore, the drivers

30    of flowering time-based reproductive isolation remain unclear. We sequenced an RNA-seq

31    dataset that comprehensively sampled *Howea* and their closest mainland relatives

32    (*Linospadix*, *Laccospadix*), and collected detailed soil chemistry data on Lord Howe Island to

33    evaluate whether secondary gene flow had taken place and to examine the role of soil

34    preference in speciation. *D*-statistics analyses strongly support a scenario whereby ancestral

35    *Howea* hybridised frequently with its mainland relatives, but this only occurred prior to

36    speciation. Expression analysis, population genetic and phylogenetic tests of selection,

37    identified several flowering time genes with evidence of adaptive divergence between the

38    *Howea* species. We found expression plasticity in flowering time genes in response to soil

39    chemistry as well as adaptive expression and sequence divergence in genes pleiotropically

40    linked to soil adaptation and flowering time. Ancestral hybridisation may have provided the

41    genetic diversity that promoted their subsequent adaptive divergence and speciation, a

42    process that may be common for rapid ecological speciation.

43

44

## Introduction

The geographic context of speciation has been a controversial topic in evolutionary biology. Theoretical models suggest that speciation can occur in sympatry with initial gene flow (Dieckmann and Doebeli 1999; Kondrashov and Kondrashov 1999; Doebeli et al. 2005; Bolnick and Fitzpatrick 2007), however it is likely to require far stronger divergent selection than speciation with spatial separation, and convincing examples in nature have been rare. It is unclear whether this reflects its genuine rarity, or the difficulty of convincing demonstration. Coyne and Orr (2004) set out four criteria for identifying cases of sympatric speciation: species are 1) currently sympatric, 2) sister taxa and 3) reproductively isolated and that 4) a period of allopatry during divergence is highly unlikely. The fourth criterion is the most difficult to demonstrate, particularly in species with broad continental distributions. Therefore, tests of sympatric speciation have attempted to reduce the possibility of an allopatric phase by focusing on species pairs restricted to small and isolated habitats such as islands. The assumption here is that the habitat island is too small for geographical isolation to have occurred within it, and too distant from other suitable habitats for speciation to have occurred during geographic isolation between the current habitat and a second one. This approach has thus far been used to infer sympatric speciation in plants and finches on remote oceanic islands (Ryan et al. 2007; Papadopulos et al. 2013) and cichlid fishes inhabiting crater lakes (Schliewen et al. 1994; Barluenga et al. 2006; Malinsky et al. 2015).

One of the best-known examples of speciation in sympatry in plants are the only two extant species of *Howea* palms (Savolainen et al. 2006). *Howea belmoreana* and *H. forsteriana* are restricted to the small (ca. 15 km$^2$) and remote (600 km from the nearest other landmass) Lord Howe Island (LHI; Australia), where they overlap in distribution across two soil types; although the former species is restricted to volcanic soil, the latter is present on both volcanic and calcareous soil. The evidence of sympatric speciation is further upheld by a sister relationship between the species in molecular phylogenetic trees (Savolainen et al. 2006; Baker et al. 2011), as well as prezygotic reproductive isolation by differences in flowering phenology (Savolainen et al. 2006; Hipperson et al. 2016) and some evidence of postzygotic isolation in the form of reduced hybrid fitness (Hipperson et al. 2016). They are both diploid, so polyploid speciation is excluded (Savolainen et al. 2006). However, a rigorous demonstration of a sister-relationship is required to provide evidence for sympatric speciation, since secondary contact (i.e. introgression) by more distantly related species can be concealed in phylogenetic trees constructed by a small number of genetic markers, as has been recently shown for seven crater lake cichlid radiations (Malinsky et al. 2015; C.H. Martin et al. 2015; Kautt et al. 2016; Meier et al. 2017; Poelstra et al. 2018). Relative to the

3

81   closest extant outgroup species to *Howea*, the monotypic *Laccospadix australasicus* and
82   species of *Linospadix*, the monophyly of *Howea* is supported by phylogenetic reconstruction
83   using only two nuclear markers. In comparison to the cichlid examples above, two markers
84   provide insufficient evidence to detect introgression. Furthermore, the existing phylogenetic
85   tree includes only three of the seven species of *Linospadix*: *L. albertisianus, L. minor* and *L.*
86   *palmerianus* (Savolainen et al. 2006). *Linospadix monostachyos,* which was not included in
87   this phylogenetic analysis, inhabits the most proximal part of Australian mainland to LHI,
88   making it the most likely known candidate for gene flow with *Howea*. A phylogenomic
89   analysis to explicitly test secondary introgression into *Howea* from mainland relatives is
90   therefore needed.

91

92   Beyond the prevalence of speciation in sympatry, many questions remain regarding its
93   genomic underpinning and the conditions that may promote it. For example, the role of
94   ancestral gene flow in providing genetic diversity on which selection can act may be more
95   important than previously appreciated (Meier et al. 2017), and the mechanisms by which
96   reproductive isolation evolves following initial local adaptation remain opaque in most
97   species. These factors are largely unknown in *Howea*, but since the species have
98   overlapping, yet distinct soil preferences, it is possible that soil characteristics have in part
99   been responsible for the divergence of species. One hypothesis is that soil acted as a driver
100  of speciation in *Howea* via selection on genes pleiotropically affecting both flowering time
101  and soil preference, thereby producing reproductive isolation as a by-product of soil
102  adaptation (Dunning et al. 2016). Another hypothesis is that a switch in soil preference
103  induced a plastic response in flowering time, allowing initial divergence. These differences
104  could have later been canalised (Pfennig et al. 2010) or bolstered by a reinforcement-like
105  mechanism (Kirkpatrick 2001). Several genes have been identified with evidence of
106  divergent selection between the *Howea* species (Dunning et al. 2016), although a lack of
107  corresponding outgroup data has prevented our ability to pinpoint in which species changes
108  have occurred. Furthermore, the characteristics of the soil have not been analysed beyond
109  broad categories (calcareous versus volcanic) and pH measurements (Savolainen et al.
110  2006; Papadopulos et al. 2013), lacking variation in soil components (e.g. macronutrients,
111  micronutrients and water availability), which may have been crucial to local adaptation and
112  the evolution of flowering time differences.
113
114  To further develop our understanding of speciation in sympatry, and test the implicated
115  mechanisms driving this process in *Howea*, we combined a detailed chemical analysis of soil
116  on LHI with an RNA-seq dataset (transcriptomes) of three tissue types derived from 54

117 individuals, which included all species of *Linospadix* and *Laccospadix* as well as both *Howea*
118 species. Specifically, we aimed to (i) determine whether external secondary gene flow from
119 mainland relatives caused speciation in *Howea*; (ii) determine whether gene flow between
120 the ancestor of *Howea* and its mainland relatives preceded colonisation of LHI; (iii)
121 characterise the soil types that the two *Howea* species inhabit, and correlate this to gene
122 expression and sequence divergence in the species; and (iv) identify which genes have
123 undergone adaptive evolution in each of the two *Howea* species, and evaluate whether
124 these include loci that could have driven the evolution of reproductive isolation.
125
126 **Materials and methods**
127
128 **Tissue sampling and RNA sequencing**
129
130 For *H. belmoreana* and *H. forsteriana*, we took the data from Dunning et al. (2016), that is,
131 19 and 17 individuals from each species, respectively. For outgroup species, we sampled
132 between one and six individuals from the wild and at the Royal Botanic Gardens, Sydney
133 (*Linospadix* albertisianus: 1, *L. apetiolatus*: 2, *L. microcaryus*: 1, *L. minor*: 3, *L.*
134 *monostachyos*: 6, *L. palmerianus*: 2, *Laccospadix australasicus*: 4; Table S1). Leaf,
135 inflorescence and root tissue was taken for each individual for RNA-sequencing. Tissue was
136 cut into <5mm$^2$ sections and stored in RNAlater (Sigma) at -20°C. All tissue samples were
137 sent to the BGI Tech solutions (Hong Kong) for RNA extraction, library construction and
138 sequencing. Paired-end 100 base pairs (bp) libraries were multiplexed and sequenced on an
139 Illumina HiSeq 4000 sequencer. These data were supplemented with publicly-available
140 RNA-seq data for other palms (Arecaceae). All data are available from SRA accession no
141 (PRJNA528594).
142
143 **Bioinformatic processing**
144
145 Illumina primer and adaptor sequences were removed and initial quality control was
146 completed by removing reads with an average PHRED-scaled quality score < 20 (both
147 completed by BGI Tech Solutions). Reads were then corrected for each individual using
148 Rcorrector v.1.0.2 (Song and Florea 2015), followed by trimming in Trimmomatic v.0.33
149 (Bolger et al. 2014) with the following settings LEADING:5, TRAILING:5,
150 SLIDINGWINDOW:4:5, MINLEN:75 (following recommendations by MacManes (2014)).
151
152 To improve the reference transcriptome of Dunning et al. (2016), we used a comprehensive
153 multi-assembler, multi-K-mer approach as follows. For the individual with the highest number

154 of corrected, trimmed read pairs for each of the two *Howea* species, we separately

155 assembled de novo transcriptomes using eight transcriptome-specific assemblers:

156 BinPacker v.1.0 (Liu et al. 2016), Bridger v. 2014-12-01 (Chang et al. 2015), IDBA-tran

157 v.1.1.0 (Peng et al. 2013), Oases v.0.2.08 (Schulz et al. 2012), Shannon v.0.0.2 (Kannan et

158 al. 2016), SOAPdenovo-Trans v.1.0.4 (Xie et al. 2014), TransABySS v.1.5.5 (Robertson et

159 al. 2010) and Trinity v.2.4.0 (Grabherr et al. 2013). For Trinity, Bridger and BinPacker, we

160 used three K-mer lengths, representing the maximum, minimum and default settings of 19,

161 25 and 33. For IDBA-tran, Oases, Shannon, SOAPdenovo-Trans and TransABySS, we used

162 six K-mer lengths, 21,31,41,51,61 and 71.  To determine fragment length distributions,

163 necessary for Bridger, Oases and BinPacker, we mapped all reads to their 25 K-mer Trinity

164 assembly using BWA-MEM v.0.7.8 (Li and Durbin 2009) with default settings. Input sizes

165 were determined from the resulting mappings using the *CollectInsertSizeMetrics* function of

166 Picard Tools v.2.6.0 (available at http://broadinstitute.github.io/picard/). This first-pass

167 assembly produced 34 independent assemblies for each species (three for Trinity, Bridger

168 and BinPacker, six for Velvet-Oases, transABYSS, SOAPdenovo and Shannon, and one for

169 IDBA-tran, which builds on each K-mer length assembly iteratively to produce one final

170 assembly).

171

172 For each of the 68 de novo assemblies, we then used TransDecoder v.3.0.1 (Haas et al.

173 2013) to identify open reading frames (ORF) using the *single_best_ORF* option. Contigs

174 lacking an ORF over 100 amino acids in length were discarded, since these are likely to

175 represent assembly artefacts. The retained coding sequences (CDS) for all 34 assemblies

176 per species were then combined and clustered using CD-HIT-EST v.4.6.1 (Fu et al. 2012)

177 with a local sequence identity threshold of 0.99, and coverage length settings of aL= 0.005,

178 aS = 1 as recommended by Cerveau and Jackson (2016). Only the longest sequence from

179 each cluster was retained to remove redundancy. To reduce the chance of assembly errors

180 in our final assembly, we removed sequences that were only recovered by a single

181 assembler or in less than four individual (i.e. assembler and K-mer length-specific)

182 assemblies, following Cerveau and Jackson (2016). This resulted in a single non-redundant

183 reference dataset for each *Howea* species. These were then matched between species

184 using reciprocal best BLAST. Each assembly was blast searched against the other with

185 BLASTN (Camacho et al. 2009) and those that were reciprocally each other's top hit were

186 retained as reciprocal best BLAST pairs (RBB-pairs). RBB-pairs were aligned using MAFFT

187 v.7.245 (Katoh et al. 2002) using automatic selection of the appropriate alignment strategy.

188 Consensus sequences for each pairwise alignment were produced using an in-house python

189 script (available at https://github.com/ogosborne/fasta_alignment_filters/), in which all non-

190 matching nucleotides were coded as ambiguity characters. To control for insertions and

191 deletions within coding regions, we reran TransDecoder v.3.0.1 (Haas et al. 2013) on the
192 consensus sequences as above. The resulting consensus sequences were used as a
193 mapping reference. To determine locus-to-transcript relationships, we used a mapping-
194 based sequence clustering pipeline. Mapping of all reads from *Howea, Linospadix* and
195 *Laccospadix* to the reference was first conducted using STAR v.2.5.3a (Dobin et al. 2013),
196 with each tissue from each individual mapped separately, allowing unlimited matches per
197 read. The mappings were then used to cluster transcripts with Corset v.1.06 (Davidson and
198 Oshlack 2014) using a distance threshold of 0.3 and considering each species-tissue type
199 combination as a distinct experimental grouping. For each resulting Corset cluster, only the
200 transcript with the longest ORF was retained. The resulting sequences formed our final
201 transcriptome assembly, and these were carried forward into downstream analyses. These
202 were matched to the previous *Howea* transcriptome assembly (Dunning et al. 2016) using
203 reciprocal best BLASTn (Camacho et al. 2009), and both assemblies were assessed using
204 BUSCO v.2.0 (Simao et al. 2015).
205
206 To identify sequence variation, STAR was run on all data, including the publicly available
207 Arecaceae data from the NCBI Short Sequence Archive (SRA), mapping all reads from each
208 individual together and allowing only unique read mappings. Variants were then called and
209 filtered using the samtools-bcftools v.1.3.1 pipeline (Li 2011). First samtools *mpileup* function
210 was used to create a pileup file for each individual separately, considering only reads with a
211 PHRED scaled mapping quality over 20. Pileups were then used to call SNPs using bcftools
212 *call* function with the multiallelic caller model, keeping sites that are ambiguous in the
213 reference, and outputting gVCF homozygous reference blocks. The resulting SNPs were
214 filtered using bcftools *filter* function excluding SNPs within three bases of an indel, with a
215 genotype quality < 20 or with a base quality < 20. Homozygous calls were required to be
216 supported by three reads and heterozygous calls were required to be supported by two
217 reads for each allele. FASTA formatted sequences were produced from the resulting VCF
218 files and reference sequences using vcf2fas, where heterozygous bases were coded as
219 IUPAC ambiguity codes (Bruno Nevado, available from
220 https://github.com/brunonevado/vcf2fas). Resulting sequences were then concatenated to
221 produce an aligned FASTA file for each gene.
222
223 Prior to phylogenetic analysis, alignments were filtered using inhouse python scripts
224 (available at https://github.com/ogosborne/fasta_alignment_filters/). To produce input data
225 for gene tree-based analyses (which used *Howea*, *Linospadix* and *Laccospadix*, as well as
226 the closest outgroup, *Areca catechu*, in order to root the trees), we first removed sequences
227 where over 2% of bases were heterozygous, as these may represent erroneous mapping of

228   paralogues to the same reference. We then removed alignment columns with over 90%

229   missing data and sequences with over 50% missing data. Following these filters, alignments

230   were retained if they contained four or more sequences and were over 100 bases in length.

231   We refer to this as the 'individual sequence dataset'. To produce input data for species tree-

232   based analysis (using all Arecaceae species), we first removed highly heterozygous

233   sequences as above, before producing a consensus sequence for each species where any

234   intraspecific variants were coded as missing data. The resulting sequences were

235   concatenated into an alignment for each gene that was then filtered for missing data as

236   above, we refer to this as the 'species-consensus sequence dataset'.

237

238   **Functional annotation**

239

240   All genes in the final reference transcriptome were annotated by BLAST. Model species

241   proteomes (primary isoforms only) and accompanying Gene Ontology (GO) terms were

242   downloaded from Phytozome v. (Goodstein et al. 2012) (Available from

243   https://phytozome.jgi.doe.gov, downloaded 11/08/18). Amino acid sequences of each gene

244   were searched against four model species proteomes, *Arabidopsis thaliana*, *Brachypodium*

245   *distachyon*, *Oryza sativa* and *Zea mays,* using BLASTP v2.2.25 (Camacho et al. 2009) with

246   an e-value cut-off of 0.0001, and only the top hit was retained. Genes were annotated with

247   the GO terms of their homologues from each of the reference proteomes. Redundancy was

248   removed and ancestor terms were added. GO term enrichment amongst genes of interest

249   was then tested using the topGO v.2.26.0 package in R (Alexa et al. 2006) using the *weight*

250   algorithm.

251

252   To specifically identify genes with the potential to drive reproductive isolation by

253   pleiotropically linking soil adaptation and flowering time in *Howea*, we identified all genes

254   with known involvement in flowering time and LHI-relevant soil characteristics. To

255   computationally identify flowering-time related genes, we took two approaches. Firstly, we

256   identified all genes for which the *A. thaliana* homologue was in the FLOR-ID flowering time

257   gene database (Bouché et al. 2016). Secondly, we identified all genes annotated with the

258   GO terms GO:0009909: "regulation of flower development" and GO:0010228: "vegetative to

259   reproductive phase transition of meristem". To computationally identify genes potentially

260   involved in soil adaptation in *Howea*, we took three approaches. Firstly, we identified all

261   genes that had significantly differential expression with regard to soil chemistry in *H.*

262   *forsteriana* (see "Analysis of gene expression" above) that we considered to have direct

263   evidence of a link with soil chemistry in *Howea*. Secondly, since water content differed

264   significantly between the soil types of LHI (see results), we identified all genes whose

265  homologues were in the DroughtDB drought gene database (Alter et al. 2015). Thirdly, we

266  identified all genes annotated with the following GO terms, which were relevant to our

267  findings from the soil chemistry analysis: GO:0006970: "response to osmotic stress",

268  GO:0009414: "response to water deprivation", GO:0042221: "response to chemical",

269  GO:0036377: "arbuscular mycorrhizal association", GO:0031667: "response to nutrient

270  levels", and GO:0006811: "ion transport". All genes which were annotated as potentially

271  involved in both flowering and soil related functions were then individually assessed with an

272  extensive literature search. We only considered genes to be potential pleiotropic for soil

273  adaptation and reproductive isolation when they had (i) published evidence of a mutant

274  flowering time phenotype and (ii) either differential expression according to soil in *Howea* or

275  published evidence of a mutant phenotype relevant to the differences we found between LHI

276  soil types.

277

278

279  **Phylogenetic inference**

280

281  Firstly, we used the species-consensus dataset to produce a dated species tree. We inferred

282  a maximum likelihood gene tree with this data using RAxML v.8.2.9 (Stamatakis 2014) with

283  200 bootstraps using the GTRGAMMA model of evolution. The branch lengths were then re-

284  estimated on 100 bootstraps of the data using the GTRGAMMA model in RAxML with the

285  topology fixed to that of the best maximum likelihood tree. These bootstrapped trees were

286  then rooted with the clade formed by *Daemonorops jenkinsiana* and *Mauritia flexuosa* as the

287  outgroup, as in previous studies (Baker et al. 2011; Couvreur et al. 2011; Faurby et al.

288  2016). Divergence times were estimated from these trees using the penalised likelihood

289  method implemented in r8s v. 1.80 (Sanderson 2003). We used three fossil calibrations,

290  taken from Faurby et al. (2016): minimum ages of 65 mya for the most recent common

291  ancestor (MRCA) of *Daemonorops* and *Mauritia*; 54.8 mya for the MRCA of *Cocos* and

292  *Elaeis*; and 85.8 mya for the MRCA of *Phoenix* and *Borassus*. We also set the root age to

293  100 mya, the crown age of Arecaceae found in previous work (Couvreur et al. 2011). For

294  each tree, we identified the optimal rate-smoothing parameter using cross-validation in r8s

295  (with the following settings: *method* = pl, *penalty* = add, *algorithm* = tn, *cvstart* = -8, *cvinc* =

296  0.5, *cvnum* = 32). The optimal rate-smoothing parameter for each tree was then used to

297  estimate divergence times and the solutions were checked with the *checkGradient* function.

298  Mean node age estimates from the 100 bootstrap replicates were taken as point estimates

299  and standard deviations were used to compute 95% confidence intervals. Because

300  topological discordance among gene trees can affect branch length estimation, we produced

301  a second species tree for which we attempted to limit its influence by removing highly

302  discordant trees from the analysis. For each gene, the best maximum likelihood gene tree
303  was compared to the best maximum likelihood species tree (see above) in a Shimodaira-
304  Hasegawa (SH) test (Shimodaira and Hasegawa 1999) implemented in RAxML (Stamatakis
305  2014). Genes for which the species tree (estimated with all genes) had a significantly worse
306  likelihood than the gene tree estimated with only the gene in question ($P < 0.05$) were then
307  removed from the analysis, and RAxML and r8s were rerun on this filtered dataset as above.
308
309  Secondly, we used the individual sequence dataset to infer a multi-species coalescence-
310  based tree, because gene tree discordance can obscure phylogenetic inference in closely
311  related species. For each gene in the individual sequence dataset, we inferred a maximum
312  likelihood gene tree with RAxML v.8.2.9 (Stamatakis 2014) with 200 bootstraps using the
313  GTRGAMMA model of evolution. SH-like branch supports were also calculated using
314  RAxML (Anisimova et al. 2011). To examine gene tree discordance using DensiTree plots,
315  all trees that contained every individual were filtered to remove trees with fewer than 10
316  nodes with under 80% SH-like support. These were then rooted using *Areca catechu* as the
317  outgroup and made ultrametric using the *root* and *chronos* functions in the APE package
318  (Paradis et al. 2004) in R v.3.3.1 (R Core Development Team 2008). These trees were then
319  visualised in DensiTree v.2.2 (Bouckaert 2010). To infer the species phylogeny while
320  explicitly accounting for gene-tree discordance we produced a coalescent-based tree. We
321  collapsed low support nodes (SH-like support < 80) as recommended by Zhang et al. (2018)
322  and retained all genes with over 50% of nodes un-collapsed. These were then input into
323  ASTRAL v.5.5.6 (Mirarab and Warnow 2015) for phylogenetic reconstruction with enforced
324  intraspecific monophyly, which inferred the topology and calculated concordance factor and
325  posterior probability based branch-support.
326
327  **Detection of introgression**
328
329  Firstly, to detect introgression we used a multidimensional scaling (MDS) approach with the
330  individual sequence dataset. SNPs for all genes were filtered to remove singletons and
331  those with over 90% missing data across all individuals. To produce a set of unlinked SNPs,
332  these were then further filtered to keep only the SNP with the least missing data per gene.
333  MDS analysis was then carried out using the *mds-plot* function in PLINK v1.9 with two
334  dimensions. If one *Howea* species were the result of hybridisation with one of the outgroups
335  we would expect it to cluster more closely with outgroups than the other *Howea* species in
336  this analysis.
337

338    Secondly, we used a *D*-statistic approach with the species-consensus dataset with the aim

339    of differentiating between three scenarios: (i) sympatric speciation in *Howea* following

340    allopatric separation from their sister taxa with no subsequent gene flow between *Howea*

341    and outgroups; (ii) sympatric speciation in *Howea* following ancestral gene flow with their

342    sister taxa; and (iii) speciation in *Howea* being driven by introgression from outgroups (Fig.

343    1). We calculated Patterson's *D* statistic (Green et al. 2010) for each four taxon subtree of

344    species with the topology (((*H. belmoreana*, *H. forsteriana*),$P_3$), outgroup) where each

345    species from *Linospadix* and *Laccospadix* was used as the third 'population', $P_3$, separately.

346    Counts of two discordant site patterns *ABBA* (((A,B),B),A) and *BABA* (((B,A),B),A) were

347    compared using Patterson's *D* statistic, with a value of *D* significantly different from zero

348    implying introgression between one of the *Howea* species and $P_3$, i.e. supporting scenario

349    (iii) above. To test for more complex introgression scenarios, including introgression

350    between the ancestor of *Howea* and the outgroups, we used five-taxon $D_{FOIL}$ statistics

351    (Pease and Hahn 2015) for each five-taxon subtree with the topology: (((*H. belmoreana*, *H.*

352    *forsteriana*),($P_3$,$P_4$)),outgroup) in which the split of $P_3$ and $P_4$ predates the split of *Howea*.

353    Four $D_{FOIL}$ statistics were calculated: $D_{FO}$, $D_{IL}$, $D_{FI}$ and $D_{OL}$ (Pease and Hahn 2015). The

354    combination of positive, negative and zero results of these four statistics can be used to infer

355    ancestral introgression. Specifically, test results in which $D_{FO}$ and $D_{IL}$ were either both

356    positive or both negative, while $D_{FI}$ and $D_{OL}$ were both zero, imply introgression between $P_3$

357    or $P_4$ and the ancestor of *Howea* (scenario ii) (Pease and Hahn 2015). Zero values for

358    Patterson's *D* are also expected under scenario (ii). Zero values for all *D*-statistics would

359    support scenario (i) (Fig. 1). *Areca catechu* was used as the outgroup in four and five-taxon

360    tests because it was the closest relative to *Howea*, *Linospadix* and *Laccospadix* in our

361    dataset. Site patterns were counted using the *fasta2dfoil.py* script from dfoil (Pease and

362    Hahn 2015). To estimate confidence intervals and *P*-values for estimates of the *D*-statistics,

363    1,000 bootstrap replicates were used, in which genes were sampled with replacement to the

364    total number of genes, and *D* was re-estimated. *P*-values were calculated for each of these

365    tests and corrected for multiple testing using Bonferroni correction. Any test with a corrected

366    *P* < 0.05 would be considered to show evidence of introgression.

367

368

369    **Soil analysis**

370

371    To analyse the soil characteristics in which *Howea* grows, we collected soil samples from 34

372    sites from which the *Howea* individuals with sequence data in this study were sourced. Soil

373    samples were sent to the Diagnostic and Analytical Services Environmental Laboratory

374    (Wollongbar, NSW, Australia) for analysis. This included three analyses: (i) 20 acid

375  extractable elements (Al, As, B, Ca, Cd, Co, Cr, Cu, Fe, K, Mg, Mn, Mo, Na, Ni, P, Pb, S, Se

376  and Zn) were quantified using inductively coupled plasma atomic emission spectroscopy

377  (ICP-AES); (ii) four Diethylenetriamine pentaacetate extractable micronutrients (Cu, Fe, Mn

378  and Zn) were quantified using ICP-AES, in a protocol which represent a closer

379  approximation of the phytoavailability of these elements; (iii) soil electrical conductance, a

380  metric that correlates with multiple soil properties related to plant health (Peverill et al. 1999),

381  was measured. For the majority of sites, water availability (27 sites) and pH (33 sites) was

382  also measured. A total of 50ml of soil was collected and weighed; samples were then dried

383  in an oven for 48 hours at 80°C. They were then re-weighed and the percent water content

384  was calculated. All samples for water content analysis were collected during a two-week

385  period following at least two weeks without rainfall from the 2$^{nd}$ to 15$^{th}$ of April 2018. Soil pH

386  was measured using Inoculo soil pH test kits (EnviroEquip Pty). To investigate overall soil

387  variation, we used a principal component analysis (PCA) approach. Missing values for pH

388  and water content were first converted to their respective median values. PCA was then

389  performed using the *prcomp* function in R with scaling to account for variable scales for

390  different components of soil variation. Each component of soil variation was also compared

391  separately. Three comparisons were applied: (i) calcareous versus volcanic soil; (ii) *H.*

392  *belmoreana* presence versus absence; and (iii) *H. forsteriana* presence versus absence.

393  These were compared for each soil type using Mann-Whitney U tests and *P*-values were

394  corrected for multiple testing using the false discovery rate (Benjamini and Hochberg 1995).

395

396  **Population structure**

397

398  To estimate the extent of isolation-by-distance within each *Howea* species, and to determine

399  whether there was any evidence of population structure according to soil type within *H.*

400  *forsteriana*, we calculated the pairwise coefficient of relatedness for all individuals in PLINK

401  using the set of unlinked SNPs used for our MDS analysis. To detect isolation-by-distance,

402  pairwise matrices of log geographic distance and coefficients of relatedness were used to

403  conduct Mantel tests using *mantel.randtest* in the R package adegenet (v. 2.1.1; Jombart

404  and Ahmed 2011); this was done for each species separately. To test whether there was

405  more divergence between soil types for *H. forsteriana*, we compared coefficients of

406  relatedness between all pairs of individuals from different soils types with all pairs from the

407  same soil type using a t-test.

408

409  **Tests of selection**

410

411 We calculated several population genetic statistics to look for evidence of selection, using

412 the PopGenome package in R (Pfeifer et al. 2014) and the individual sequence dataset.

413 Differentiation (as measured by $F_{ST}$; Weir and Cockerham 1984) and net divergence ($d_{XY}$;

414 Nei, 1987) were calculated for each contig between the two *Howea* species, and Tajima's D

415 (Tajima 1989) and average pairwise diversity ($\pi$) were calculated within each species.

416 Genes with $F_{ST}$ of 1 (indicating complete fixation), $d_{XY}$ over the 95[th] percentile, or Tajima's D

417 below the 95[th] percentile were considered genes of interest for downstream analyses.

418

419 To identify genes potentially evolving under positive selection in the two *Howea* species, we

420 also used a phylogenetic $d_N/d_S$-based approach with the species-consensus dataset, which

421 was trimmed to include only CDS sequences inferred by TransDecoder (Haas et al. 2013).

422 We then implemented various codon filters. Sequences with premature stop codons were

423 removed, as well as the final stop codon of each alignment. To ensure that the amount of

424 selection could be compared between the two *Howea* species, we retained only codons that

425 contained no missing data in either species. Following these filters, alignments containing at

426 least 33 codons, at least one SNP between the two *Howea* species, and at least three

427 species in the alignment were taken forward for tests of selection. We used the branch-site

428 test of positive selection (Zhang et al. 2005) implemented in the codeml program in PAML

429 v.4.8 (Yang 2007). The branch-site models allow selection to vary across both sites and

430 branches of the phylogeny. For each *Howea* species, we implemented two models. The

431 alternative model allows $d_N/d_S$ to vary above 1 on some sites on the branch being tested (the

432 foreground branch, designated as the tips leading to each *Howea* species separately)

433 whereas other branches (background branches) only vary between 0 and 1. The alternative

434 model is compared to a null model where $d_N/d_S$ is fixed at 1 for these sites on the foreground

435 branch. It is compared in a Likelihood Ratio Tests (LRT), which approximates a chi-squared

436 distribution with one degree of freedom. Phylogenetic uncertainty can affect the results of

437 the branch-site test (Pie 2006), so we took two approaches to ensure our results were robust

438 to it. First, we reran the significant branch-site tests using a species tree estimated with only

439 genes that were not significantly discordant with the overall species topology (see

440 "Phylogenetic inference" section above). Second, we reran the significant branch-site tests

441 using the best maximum likelihood gene tree for the gene tested, rather than the species

442 tree.

443

444

445 **Analysis of gene expression**

446

447  Read counts for each tissue type and each gene for all individuals in *Howea*, *Linospadix* and
448  *Laccospadix* were produced by CORSET. These were used for phylogenetic analysis of
449  gene expression. Counts were converted to reads per million to correct for differences in
450  total numbers of reads per individual. Normalised read data were then used for LRT for
451  branch-specific expression shift tests implemented in EVE (Rohlfs and Nielsen 2015). This
452  approach models the evolution of gene expression as an Ornstein–Uhlenbeck process,
453  where the parameter $\theta_i{}^a$ represents the optimal expression level for gene *i* in lineage *a*, and
454  $\theta_i{}^{non\text{-}a}$ represents the optimal expression level for gene *i* in all other lineages. The test
455  compares a null model where $\theta_i{}^a = \theta_i{}^{non\text{-}a}$ with an alternative model where $\theta_i{}^a \neq \theta_i{}^{non\text{-}a}$. The
456  models are compared with an LRT to identify significant expression shifts in the focal lineage
457  *a*. Six LRTs were implemented, testing for significant expression shifts in the branches
458  leading to *H. belmoreana* and *H. forsteriana* separately in each of the three tissue types.
459  LRT statistics were used to calculate *P*-values using chi-squared tests with one degree of
460  freedom, and the *P*-values were corrected for multiple testing using FDR. Genes that
461  showed significant expression shifts ($P < 0.05$) following multiple test correction were
462  considered genes of interest in downstream analyses.
463
464  Finally, we also investigated whether gene expression within each species was related to
465  variation in soil chemistry. The soil chemistry dataset is highly multidimensional, so we used
466  the first principal component of the soil PCA (above), which separates the two soil types, as
467  an explanatory variable in an analysis of differential expression. We used DESeq2 (Love et
468  al. 2014) to test for differential expression across PC1 of soil chemistry within each *Howea*
469  species and tissue type combination separately. Following Dunning and colleagues (2016),
470  who published the transcriptome data for the two *Howea* species, we also included sampling
471  date (categorised into three collecting trips) as a confounding variable in the model. DESeq2
472  *P*-values were corrected for multiple testing using FDR.
473
474  **Results**
475
476  **Transcriptome**
477
478  For each of the 19 *Linospadix* and *Laccospadix* individuals (Table S1), between 26,173,535
479  and 45,115,994 paired-end fragments were sequenced using RNA-seq. These were
480  supplemented with previously published RNA-seq data for the two *Howea* species (36
481  individuals) and data from ten other Arecaceae species. Following read correction and
482  trimming, all newly sequenced individuals had between 25,990,695 and 44,902,175 reads
483  remaining (Table S2). Utilising a multi-assembler and multi-Kmer pipeline using the

484 individual of each *Howea* species with the most reads, we produced a final transcriptome
485 assembly containing 26,972 genes. BUSCO analysis found that the reference transcriptome
486 had 88.7% completeness, a substantial improvement on the 77.5% completeness of the
487 previous assembly by Dunning et al. (2016), demonstrating the utility of our approach. Read-
488 mapping to the transcriptome assembly resulted in between 65% and 84% of reads being
489 uniquely mapped for *Howea* individuals, between 58% and 81% in *Linospadix* and
490 *Laccospadix*, and between 19% and 74% for other palms (Table S3).
491
492 **Ancestral hybrid swarm followed by sympatric speciation in *Howea***
493
494 The two *Howea* species were supported as sister species with a 100% of bootstrap support
495 for every node in our RAxML tree using all concatenated transcripts (Fig. S1). The
496 coalescent-based ASTRAL tree was topologically identical, and posterior probability support
497 for all nodes was high (>0.99). Unlike previous analyses that resolved *Linospadix* as sister
498 to *Howea* (Savolainen et al. 2006), we found a sister relationship between *Laccospadix* and
499 *Linospadix*. The short branch length between the common ancestors of *Howea-Linospadix-*
500 *Laccospadix* and *Linospadix-Laccospadix* and high level of ancestral introgression (see
501 below) likely explains the difference between studies. We estimated the divergence time of
502 the two *Howea* species as 3.3 million years ago (Fig. 2a). The second dated tree we
503 produced, in which genes that were phylogenetically incongruent with the species topology
504 were removed (Fig. S2), had a divergence time for the two *Howea* species of 4.4 million
505 years ago. These dates are older than previous estimates (Savolainen et al. 2006) but are
506 still well within the age of LHI, which was formed between 6.4 and 6.9 million years ago
507 (McDougall et al. 1981).
508
509 There was no evidence of hybridisation between extant *Howea* and *Linospadix* or
510 *Laccospadix*. Firstly, tree topologies were highly consistent between loci (Fig. 2a) with 95%
511 of gene trees supporting the monophlyly of *Howea* (Fig 2a-b; i.e. 95% quartet support).
512 Quartet support for the other interspecific nodes between species in *Linospadix,*
513 *Laccospadix* and *Howea* ranged from 38% to 87%, demonstrating a relatively high level of
514 gene tree-species tree discordance within the dataset. Patterson's *D* Statistics were not
515 significantly different from zero when any *Laccospadix* or *Linospadix* species were tested as
516 a potential introgressant (Fig. 3b, Table S4). However, all five-taxon *D*-statistic tests showed
517 evidence for introgression between *Linospadix* or *Laccospadix* and the *ancestor* of *Howea*
518 (scenario (ii) in Fig. 1; Fig. 3a; Table S5). All comparisons in which *Laccospadix* and one of
519 the *Linospadix* species were $P_3$ and $P_4$ showed evidence for *Laccospadix* as the
520 introgressing taxon, suggesting that admixture between *Laccospadix* and ancestral *Howea*

521   was either stronger or more recent than it was between *Linospadix* and ancestral *Howea*

522   (Pease and Hahn 2015). There was also evidence for admixture between some *Linospadix*

523   species and ancestral *Howea* when two *Linospadix* species were assigned as $P_3$ and $P_4$.

524   The only species with no evidence of admixture was *L. albertisianus*, which is restricted to

525   New Guinea, whereas all other species tested are found on the Australian mainland.

526   Defining the exact patterns of introgression is not possible, but they indicate a complex

527   history in which admixture has occurred independently between ancestral *Howea* and

528   several lineages within *Linospadix-Laccospadix* independently (one interpretation which is

529   consistent with the results is shown in Fig. 3c).

530

531   The MDS analysis of all sequence polymorphism data was consistent with the results above.

532   It resolved both *Howea* species as distinct clusters, and outgroup species were equidistant

533   from the two *Howea* species (Fig. 4). If one species was a product of hybridisation between

534   ancestral *Howea* and an outgroup (as shown in scenario (iii) in Fig. 1), it would be expected

535   to cluster more closely to the outgroup than the non-hybrid species. Instead, the reported

536   equidistance in the MDS supports lack of hybridisation during or after speciation in *Howea*.

537   Furthermore, the fact that *Laccospadix* is closer to *Howea* than is its sister taxon *Linospadix*,

538   is in line with the higher level of admixture between *Laccospadix* and ancestral *Howea* (*D*-

539   statistic results above).

540

541   Our Mantel tests revealed no evidence for isolation by distance in either *Howea* species (*H.*

542   *belmoreana*: $P = 0.440$; *H. forsteriana*: $P = 0.588$; Fig. S3), indicating that geographically-

543   based isolation within LHI is unlikely.

544

545   Taken together, our results strongly support a scenario whereby ancestral *Howea* was part

546   of a hybrid swarm on mainland Australia, but neither hybridisation with outgroups or

547   allopatric isolation within LHI drove speciation following the colonisation of LHI by the

548   ancestral *Howea*.

549

550

551   **Soil chemistry drives expression shifts in flowering time genes in *Howea***

552

553   Detailed soil chemical analysis (Table S6) revealed substantial differences between the two

554   soil types on LHI, and between specific sites that each species inhabits. The PCA analysis

555   showed that most volcanic sites were clustered, with two outliers (Fig. 5b). In contrast to this,

556   calcareous soils were more diffusely distributed. Notably, both volcanic outliers were sites

557   only occupied by *H. forsteriana* and it is possible these represented intermediate soil types

558  (Fig. 5a). All sites inhabited by *H. belmoreana* were clustered in the PCA (Fig. 5c).

559  Conversely, *H. forsteriana*-inhabited sites were widely distributed across both of the first two

560  principal components (Fig. 5d).

561

562  When individual constituents of soil variation were compared, 21 constituents were

563  significantly different between volcanic and calcareous soils, 15 were significantly different

564  between *H. belmoreana* present versus absent sites, and three were significantly different

565  between *H. forsteriana* present versus absent sites (Fig. S4). Calcareous soil was

566  characterised by significantly higher concentrations of arsenic, boron, calcium, cadmium,

567  sodium, phosphorus and sulphur, high pH and lower water content. Volcanic soils were

568  characterised by significantly higher concentrations of aluminium, cobalt, chromium, copper,

569  iron, potassium, manganese, nickel and zinc, higher water content and neutral pH (Fig. S4).

570  Overall, our soil analysis emphasises that, whereas *H. belmoreana* is an edaphic specialist,

571  *H. forsteriana* is a generalist able to grow on a far broader range of soil types.

572

573  When the first principal component of soil chemistry variation was used as an explanatory

574  variable in a differential expression analysis (for each *Howea* species and tissue type

575  separately), very few genes were differentially expressed according to soil variation in *H.*

576  *belmoreana* (inflorescence: four genes, leaf: two genes, root: seven genes). In *H.*

577  *forsteriana,* while again there was minimal soil-related differential expression in

578  inflorescences and roots (inflorescence: 15 genes, root: 22 genes), there was a very high

579  level of differential expression in leaves (1,118 genes, Table S7). This included 18 genes

580  known to be involved in flowering time differences in model plants, potentially indicating a

581  link between soil chemistry and flowering time divergence in *Howea*. We found that pairs of

582  *H. forsteriana* individuals from the same soil type were no more closely related than pairs of

583  individuals from different soil types (t-test: $P = 0.78$), indicating that there were no 'ecotypes'

584  within *H. forsteriana* that would explain gene expression differences, especially considering

585  that such a large number of genes were differentially expressed.

586

587  **Adaptive evolution of protein sequence and gene expression is species and tissue**

588  **specific in *Howea***

589

590  We used phylogenetic approaches to search for genes with amino acid substitutions under

591  positive selection (Yang 2007) as well as genes with a significant shifts in expression level

592  (Rohlfs and Nielsen 2015) in the branches leading to each *Howea* species. Genes that have

593  undergone significant expression shifts were unevenly distributed across tissue types and

594  species. In total, 1,736 genes have undergone an expression shift in at least one tissue in at

least one species. In inflorescence and root tissue, there were significantly more in *H. forsteriana* (inflorescence: 560 in *H. forsteriana* versus 100 in *H. belmoreana*, Fishers Exact Test *P* < 0.001; root: 131 in *H. forsteriana* versus 20 in *H. belmoreana*, Fishers Exact Test *P* < 0.001), whereas in leaf tissue there were significantly more in *H. belmoreana* (892 in *H. belmoreana* versus 227 in *H. forsteriana*, Fishers Exact Test *P* < 0.001). The tests of positive selection on amino acid sequence revealed that 104 genes likely evolved under positive selection in *H. belmoreana* while 132 were under positive selection in *H. forsteriana*; although this difference in numbers was not significant (*P* = 0.077; Fisher's exact test). Of the 1,972 genes that had any evidence of adaptive evolution from these tests, only 9% were found in more than one of these sets (Fig. S5).

**Genes under adaptive evolution are enriched for edaphic and phenology-related functions**

There was significant enrichment of 103 GO terms amongst our genes of interest (e.g. genes showing an expression shift or significant sequence-based evidence of positive selection, see Methods), ranging from two GO terms amongst genes which underwent a significant expression shift in the roots of *H. belmoreana* to 27 GO terms amongst genes with any evidence of adaptive evolution (Table S8). Several of these genes were relevant to the speciation scenario of *Howea.* For example, genes with either evidence of an expression shift or with a signature of positive selection in *H. forsteriana* were significantly enriched for the GO term "negative regulation of flower development", indicating that the differing flowering times of the two species has evolved adaptively (Table S8, GO term assignment and evidence of selection listed in Table S7). Several GO terms likely to be involved in soil preference differences between the two species are also enriched amongst candidate genes (Table S8). The term "response to cadmium ion" is over-represented among genes showing evidence for positive selection in *H. forsteriana*, "response to water deprivation" is over-represented amongst genes under positive selection in *H. belmoreana*, "cellular calcium ion homeostasis" is over-represented amongst genes which have undergone an expression shift in leaf tissue in *H. forsteriana*, and "cellular response to phosphate starvation is over-represented amongst genes that have undergone a significant expression shift in the inflorescence of *H. forsteriana*. Cadmium, phosphorus, calcium and water content all vary significantly between the soils of the two species. Several less specific GO terms relevant to soil chemistry differences between calcareous and volcanic soil such as "transition metal ion transport", "regulation of ion transport" and "divalent metal ion transport" were also enriched amongst genes of interest. Furthermore, several terms involved in biotic interactions known to differ between the soil types and species were over-represented. This included several

defence related GO terms: "defence response", "defence response signalling pathway", "defence response to bacterium", "response to bacterium" and "regulation of defence response". Osborne et al. (2018) showed that multiple plant pathogens, both fungal and bacterial, are differentially abundant between the two soil types and even between the two species on the same (volcanic) soil type. Our results here indicate that pathogens may act as selection pressure on the species. Finally, we found that the term "response to karrikin" was significantly over-represented in genes that had experienced an expression shift in *H. forsteriana*. Karrikin response genes are involved in the initiation of arbuscular mycorrhizal fungi symbiosis in rice (Gutjahr et al. 2015). This may be important in *Howea* too, since the two species have divergent soil-specific interactions with arbuscular mycorrhizal fungi, which likely affect their relative fitness on calcareous versus volcanic soils (Osborne et al. 2018).

We also identified 122 genes within our dataset (Tables S7 and S9) that have the potential to pleiotropically link soil adaptation and flowering time in *Howea*. This is based on their mutant phenotypes in model plant species (Table S9; references listed in table). Nine of these showed evidence of adaptive evolution in one or both of the species (Table S10). Two of these nine were amongst the six candidate 'speciation genes' identified by Dunning et al. (2016), and another one, *DCL1*, was annotated to the same *Arabidopsis* homologue as in Dunning et al., although we did not return each other as best reciprocal match in our BLAST searches. Two of the six candidates for positive selection in Dunning et al. ($d_N/d_S > 1$), however, did not show evidence of positive selection in our branch-site test, although because this test is highly conservative, this may not necessarily surprising (Gharib and Robinson-Rechavi 2013). In total, combining results from this study with that of Dunning et al (2016) identified 13 candidate 'speciation genes' in *Howea* (Table S10).

**Discussion**

**The role of admixture in *Howea***

We found no evidence for gene flow from outgroups into *Howea* following initial colonisation, supporting a model of sympatric speciation. Traditionally, the main difficulty of demonstrating sympatric speciation has been to show that there has been no potential for geographic separation within their habitat such that reproductive isolation could have evolved in allopatry. For this reason, while many potential examples of sympatric speciation may exist on continental landmasses (Sorenson et al. 2003; Hadid et al. 2013; Osborne et al. 2013; Hadid et al. 2014), the most convincing case studies have been found in tiny habitat islands such as crater lakes and oceanic islands (Schliewen et al. 1994; Savolainen et al. 2006;

669 Papadopulos et al. 2011; Malinsky et al. 2015). Recent evidence of secondary gene flow into
670 several crater lakes hosting cichlid radiations that were previously thought to have evolved in
671 complete sympatry has cast doubt on these examples (Martin et al. 2015; Poelstra et al.
672 2018). All four cichlid radiations in Cameroon were shown to involve secondary gene flow
673 from external riverine populations, with some even being more closely related to the river
674 populations than other species within their lakes (Martin et al. 2015). Furthermore, in one of
675 these lakes, a genomic region containing several olfactory genes involved in mate choice in
676 cichlids was introgressed prior to the first speciation event (Poelstra et al. 2018). This
677 indicates a mechanism by which secondary gene flow may have played a causative role in
678 their speciation (Poelstra et al. 2018). While some authors considered secondary gene flow
679 to be unlikely to have been causative in some other cichlid radiations (Malinsky et al. 2015),
680 this cannot currently be ruled out. Given that our results make a role for secondary gene flow
681 highly unlikely, *Howea* appears to be one of the strongest examples of sympatric speciation
682 in nature.
683
684 Nevertheless, admixture with outgroups may be important in the evolutionary history of
685 *Howea*. Admixture can play a key role in generating genetic diversity on which selection can
686 act (Seehausen 2004; Seehausen 2015; Arnold and Kunte 2017) and has been shown to
687 precede adaptive divergence in multiple taxa. For example in Lake Victoria cichlids,
688 ancestral admixture produced exceptional variation in opsin genes known to be involved in
689 speciation and adaptation, thereby facilitating adaptive radiation in the lake (Meier et al.
690 2017). In Darwin's finches, admixture between two species increased the standing genetic
691 and evolutionary responsiveness to fluctuating environmental conditions (Grant and Grant
692 2014). In light of these and other examples (Pardo-Diaz et al. 2012; Stankowski and
693 Streisfeld 2015), our finding that ancestral *Howea* may have been part of a mainland hybrid
694 swarm opens the possibility that ancestrally-introgressed variation could have been
695 important in their divergence following the colonisation of LHI. This may be critical given the
696 likelihood of a genetic bottleneck upon colonisation of LHI. Unfortunately, identifying
697 introgressed regions requires longer genomic windows than we can derive from
698 transcriptomic data and so it is outside the scope of this study (Martin et al. 2015), although
699 our ongoing sequencing of the *Howea* genome will provide an opportunity to test this
700 hypothesis in the future.
701
702 We have interpreted the *D*-statistic results in terms of introgression; however, non-zero
703 results could also be obtained because of ancestral population structure (Pease and Hahn
704 2015). For example, consider three ancestral populations *A, B* and *C*, in which gene flow
705 between *B* and *C* is stronger than between *A* and *C*. Then, consider that the three

populations later diverge into three species as per the phylogenetic tree ((A,B),C), with no gene flow following speciation. The *D*-statistic may then imply, wrongly, introgression between species *B* and species *C*. However, in our case, a population structure-only scenario to explain all our results would be highly complex. The relatively ancient age of the *Howea* species splits makes it implausible that all $D_{FOIL}$ results are explained by ancestral population structure alone, which has been subsequently preserved to the present day. Furthermore, since all Patterson's *D*-statistics were zero, there was no ancestral population structure by which one nascent *Howea* species experienced more gene flow with outgroups than the other. This bolsters the case that ancestral *Howea* was one homogenous population upon colonisation of LHI.

**Ball's Pyramid is an unlikely source of geographic isolation**

Ball's Pyramid is an inhospitable sea stack 30 km off LHI. It is known that Ball's Pyramid, as well as LHI, had a greater terrestrial extent in the past due to lower sea levels at some time in the Pleistocene (Papadopulos et al. 2011; Woodroffe et al. 2006; Linklater et al. 2018). This raised the possibility of allopatric speciation in *Howea* between populations isolated on LHI and Ball's Pyramid. Papadopulos and colleagues (2011) investigated this possibility and concluded that allopatric divergence of *Howea* was unlikely. The argument was that the distance that would have separated populations of *Howea* on LHI and Ball's Pyramid would not be greater than the current length of LHI, on which populations are not geographically structured in this wind-pollinated genus. Furthermore, a recent demographic modelling analysis supports a model in which gene flow following speciation was high and reduced towards the present over models which include an allopatric period (Papadopulos et al., 2019a). Our new analyses strengthen the case further. While sea level has been lower (and thus the terrestrial extent of both islands has been larger) in the past, our dates indicates that *Howea* split in the Pliocene, when sea level was actually about 22 m higher than today (Dwyer and Chandler 2009; Miller et al. 2012). Erosion of Ball's pyramid and LHI to their current state, where Ball's pyramid is a shear sea stack completely unsuitable for *Howea* colonisation, occurred rapidly following their formation and substantially earlier than the *Howea* split (6-7 mya; Linklater et al 2018). In any case, during periods of lowest sea level, the distance between the islands was around 4 km, whereas we found no evidence of isolation-by-distance within either species with a maximum distance of 5.8 km (Fig. S3). While this does not necessarily mean that pollen travels over these distances, taken together with the evidence from demographic modelling and the likely unsuitability of Ball's Pyramid for *Howea* colonisation at the time of speciation, all available evidence indicate that

742    allopatric isolation between Ball's pyramid and LHI is unlikely to have been responsible for

743    speciation in *Howea*.

744

745

746    **The ecological circumstances of speciation in *Howea***

747

748    The main ecological difference between the *Howea* species today is their soil preferences.

749    Our analyses highlight the fact that while *H. belmoreana* is a specialist on volcanic soil, *H.*

750    *forsteriana* is a soil generalist. However, the observation that *H. forsteriana* has not

751    displaced *H. belmoreana* indicates that this generalism comes at a cost. Common garden

752    experiments have provided evidence for this, showing that *H. belmoreana* has a higher

753    survival rate than *H. forsteriana* on volcanic soil (Hipperson et al. 2016). One explanation for

754    this difference is that *H. forsteriana* is less able to form arbuscular mycorrhizal associations

755    on volcanic soil than either *H. belmoreana* on volcanic soil or *H. forsteriana* on calcareous

756    soil (Osborne et al. 2018).

757

758    Several soil characteristics significantly differed depending on broad soil type categories

759    (volcanic versus calcareous) and the presence or absence of each species. The differences,

760    which are likely to exert contrasting stresses on plants, comprised changes in essential

761    primary (P, K) and secondary (S, Ca) macronutrients, micronutrients (Al, B, Co, Cu, Fe, Na,

762    Mg, Mn, Ni, Zi) and toxins (As, Cd, Cr). This highlights the fact that a switch in soil is a

763    multidimensional environmental change, and is likely to affect multiple genes, potentially

764    increasing the barriers to gene flow more than simpler environmental switches (Nosil et al.

765    2017; Riesch et al. 2017).

766

767    Given their current divergent soil preferences in the two *Howea* species, it is likely that soil

768    adaptation played a role in their divergence. Since there was admixture between ancestral

769    *Howea* and *L. minor* following the split of *L. minor* and *L. albertisianus*, we can

770    approximately date the colonisation of LHI to between 4.96 and 3.27 mya (the speciation

771    times of *L. minor* and *L. albertisianus*, and *H. belmoreana* and *H. forsteriana*, respectively;

772    this is in line with our dating calculations with r8s). The current calcareous formations on LHI

773    were deposited in the last 350,000 years (Brooke et al. 2003), substantially later than this.

774    However, we had hypothesised that it was more likely that the ancestral *Howea* first

775    colonised volcanic soils, and subsequently moved to calcareous deposits due to the fact that

776    the volcanic soils are older and more similar to mainland soils (Savolainen et al. 2006;

777    Papadopulos et al., 2019b). Since the erosion rate of calcarenite can be around 2.35

778    mm/year (Balaguer et al. 2019), calcareous deposits from c.a. 3-5 mya would unlikely have

779    survived today. Of course, it is also possible that the ancestral *Howea* colonised calcareous

780    soils first, before colonising volcanic ones. While we do not know the detailed edaphic

781    composition on the island at the time of colonisation and speciation, we can get insight into

782    the ecological history of *Howea* from the proportion and identity of genes that have

783    undergone adaptive divergence in the species.

784

785    As would be expected given the relatively old divergence of *Howea*, there are a multitude of

786    sequence and gene expression differences between the two species. Previous research in

787    *Howea* (Dunning et al. 2016) could not polarise the trajectory of change in these loci

788    because of the absence of outgroup data. Under our default scenario in which ancestral

789    *Howea* was a volcanic specialist and speciation was precipitated by an invasion of

790    calcareous soil by the ancestor of *H. forsteriana*, it may be expected that more adaptive

791    evolution should be found in this latter species. This was indeed the case for expression

792    shifts in inflorescence and root, as well as in positive selection on coding sequences (Fig.

793    S5). Furthermore, these genes in *H. forsteriana* were significantly enriched for several soil

794    adaptation-related functions (Table S8). In contrast, significantly more expression shifts were

795    found in the leaves of *H. belmoreana*, and while the reason for this is unclear, it should be

796    noted that the leaf morphology of *H. belmoreana* is unusual relative to both *H. forsteriana*

797    and outgroup species in *Linospadix* and *Laccospadix*, featuring recurved leaves with

798    ascending leaflets (Savolainen et al. 2006) (compare Fig. 2c and 2d).

799

800    We note that there are, of course, caveats to consider in gene expression results. Since our

801    samples are wild trees, tissues were collected in different locations and at different times,

802    and information such as plant age and health cannot be known. The environment affects

803    gene expression, and therefore our samples from the wild cannot be standardised as they

804    would be if derived from greenhouse experiments. Nevertheless, we have controlled for

805    sampling date in our models, and we found that the several significant shifts are consistent

806    with our hypotheses.

807

808    While our overall results are consistent with a scenario of ancestral volcanic-specialism,

809    alternatives may still be possible. What we can say conclusively, is that the evolution of the

810    two species involved adaptation to the abiotic (Hipperson et al. 2016) and biotic (Osborne et

811    al. 2018) soil variation on the island, so we then investigated a link between soil adaptation

812    and the main component of reproductive isolation, flowering time.

813

814    **The evolution of reproductive isolation in *Howea***

815

816    We identified genes that could have linked ecological adaptation to soil and reproductive

817    isolation via flowering time via two distinct mechanisms: plasticity and pleiotropy. Notably,

818    cadmium (Cd) and zinc (Zi) both differ between the soil types, with cadmium being

819    significantly higher in calcareous soil, and zinc being significantly higher in volcanic soil.

820    Substrate concentrations of these two elements have both been experimentally shown to

821    alter flowering time in *Arabidopsis* (Wang et al. 2012; Przedpelska-Wasowicz and Wasowicz

822    2013). Therefore, migration between the soil types has the potential to cause a plastic shift

823    in flowering time, a mechanism which may be common in plants (Levin 2009). Such a shift

824    could have instantaneously reduced gene flow between the two nascent *Howea* species in

825    the early stages of their evolution. This scenario is supported in *Howea* given that we found

826    multiple known 'flowering time genes', which were differentially expressed according to soil

827    chemistry in *H. forsteriana*. While mean flowering time is not significantly different between

828    the two soil types in *H. forsteriana*, it flowers protoandrously (male flowers are produced

829    earlier) on calcareous soil, whereas on volcanic soil male and female flowering have been

830    found to be synchronous in at least one population (Savolainen et al. 2006), showing that

831    soil chemistry may affect aspects of flowering phenology. Furthermore, there are many

832    examples of a loss of plasticity during or following speciation (Aubret and Shine 2009;

833    Pfennig et al. 2010; Palmer 2014; Levis et al. 2018), so soil-specific flowering time

834    differences could have been more pronounced in the past. Flowering time plasticity would be

835    expected to drive speciation most strongly when pollen dispersal is high but seed dispersal

836    is low, since pollen-mediated gene flow should be directly affected by flowering time

837    whereas seed-mediated gene flow should not. This is likely the situation in *Howea*, which is

838    wind-pollinated but has large and immobile seeds (Savolainen et al. 2006). Therefore soil-

839    mediated flowering time plasticity is a plausible mechanism of speciation in *Howea*. With the

840    combined results of this study and Dunning et al (2016), we also identified 13 genes

841    showing evidence of adaptive expression or sequence divergence in *Howea,* with functions

842    that could pleiotropically link soil adaptation and flowering time (Table S10). While we

843    tentatively consider these genes to be candidate 'speciation genes', their potential function

844    was inferred from sequence similarity to model plant genes of known function, which does

845    not guarantee that the orthologous palm genes have the same function. Given the long

846    generation times of palm trees, functional assays using palm mutants are not practical.

847    However, assays involving knockout mutants in model systems for candidate orthologous

848    genes, followed by phenotype rescue using palm genes would provide clearer insight into

849    the function of key genes. This work is ongoing by Savolainen and Turnbull and should

850    further elucidate the genomic basis for speciation in *Howea*.

851

852

**References**

Alexa A, Rahnenführer J, Lengauer T. 2006. Improved scoring of functional groups from gene expression data by decorrelating GO graph structure. *Bioinformatics* 22:1600–1607.

Alter S, Bader KC, Spannagl M, Wang Y, Bauer E, Schön CC, Mayer KFX. 2015. DroughtDB: An expert-curated compilation of plant drought stress genes and their homologs in nine species. *Database* 2015:1–7.

Anisimova M, Gil M, Dufayard JF, Dessimoz C, Gascuel O. 2011. Survey of branch support methods demonstrates accuracy, power, and robustness of fast likelihood-based approximation schemes. *Syst. Biol.* 60:685–699.

Arnold ML, Kunte K. 2017. Adaptive Genetic Exchange: A Tangled History of Admixture and Evolutionary Innovation. *Trends Ecol. Evol.* 32:601–611.

Aubret F, Shine R. 2009. Genetic Assimilation and the Postcolonization Erosion of Phenotypic Plasticity in Island Tiger Snakes. *Curr. Biol.* 19:1932–1936.

Baker WJ, Norup M V., Clarkson JJ, Couvreur TLP, Dowe JL, Lewis CE, Pintaud JC, Savolainen V, Wilmot T, Chase MW. 2011. Phylogenetic relationships among arecoid palms (Arecaceae: Arecoideae). *Ann. Bot.* 108:1417–1432.

Balaguer P, Pons GX, Mir-Gual M. 2019. The Rocky Coasts of Balearic Islands: Dynamic Processes, Sediments and Management. In: Morales JA, editor. The Spanish Coastal Systems. Springer International. p. 116–141.

889    Barluenga M, Stölting KN, Salzburger W, Muschick M, Meyer A. 2006. Sympatric speciation
890        in Nicaraguan crater lake cichlid fish. *Nature* 439:719–723.

891    Benjamini Y, Hochberg Y. 1995. Controlling the False Discovery Rate: a Practical and
892        Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B* 57:289–300.

893    Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: A flexible trimmer for Illumina sequence
894        data. *Bioinformatics* 30:2114–2120.

895    Bolnick DI, Fitzpatrick BM. 2007. Sympatric Speciation: Models and Empirical Evidence.
896        *Annu. Rev. Ecol. Evol. Syst.* 38:459–487.

897    Bouché F, Lobet G, Tocquin P, Périlleux C. 2016. FLOR-ID: An interactive database of
898        flowering-time gene networks in *Arabidopsis thaliana*. *Nucleic Acids Res.* 44:1167–
899        1171.

900    Bouckaert RR. 2010. DensiTree: Making sense of sets of phylogenetic trees. *Bioinformatics*
901        26:1372–1373.

902    Brooke BP, Woodroffe CD, Murray-Wallace C V., Heijnis H, Jones BG. 2003. Quaternary
903        calcarenite stratigraphy on Lord Howe Island, southwestern Pacific Ocean and the
904        record of coastal carbonate deposition. *Quat. Sci. Rev.* 22:859–880.

905    Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.
906        BLAST+: Architecture and applications. *BMC Bioinformatics* 10:421.

907    Cerveau N, Jackson DJ. 2016. Combining independent de novo assemblies optimizes the
908        coding transcriptome for nonconventional model eukaryotic organisms. *BMC*
909        *Bioinformatics* 17:525.

910    Chang Z, Li G, Liu J, Zhang Y, Ashby C, Liu D, Cramer CL, Huang X. 2015. Bridger: a new
911        framework for de novo transcriptome assembly using RNA-seq data. *Genome Biol.*
912        16:30.

913    Couvreur TLP, Forest F, Baker WJ. 2011. Origin and global diversification patterns of
914        tropical rain forests: Inferences from a complete genus-level phylogeny of palms. *BMC*
915        *Biol.* 9:44.

916    Coyne JA, Orr HA. 2004. Speciation. Sunderland (MA): Sinauer Associates.

917    Davidson NM, Oshlack A. 2014. Corset: Enabling differential gene expression analysis for
918        de novo assembled transcriptomes. *Genome Biol.* 15:410.

919    Dieckmann U, Doebeli MO. 1999. On the origin of species by sympatric speciation. *Nature*
920        400:354–357.

921    Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M,
922        Gingeras TR. 2013. STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics* 29:15–
923        21.

924    Doebeli M, Dieckmann U, Metz JA, Tautz D. 2005. What we have also learned: adaptive
925        speciation is theoretically plausible. *Evolution* 59:691–699.

926 Dunning LT, Hipperson H, Baker WJ, Butlin RK, Devaux C, Hutton I, Igea J, Papadopulos
927    AST, Quan X, Smadja CM, et al. 2016. Ecological speciation in sympatric palms: 1.
928    Gene expression, selection and pleiotropy. *J. Evol. Biol*. 29:1472–1487.

929 Dwyer GS, Chandler MA. 2009. Mid-Pliocene sea level and continental ice volume based on
930    coupled benthic Mg/Ca palaeotemperatures and oxygen isotopes. *Phil. Trans. R. Soc. A*
931    367:157-168.

932 Faurby S, Eiserhardt WL, Baker WJ, Svenning JC. 2016. An all-evidence species-level
933    supertree for the palms (Arecaceae). *Mol. Phylogenet. Evol.* 100:57–69.

934 Fu L, Niu B, Zhu Z, Wu S, Li W. 2012. CD-HIT: Accelerated for clustering the next-
935    generation sequencing data. Bioinformatics 28:3150–3152.

936 Gharib WH, Robinson-Rechavi M. 2013. The branch-site test of positive selection is
937    surprisingly robust but lacks power under synonymous substitution saturation and
938    variation in GC. *Mol. Biol. Evol*. 30:1675–1686.

939 Goodstein DM, Shu S, Howson R, Neupane R, Hayes RD, Fazo J, Mitros T, Dirks W,
940    Hellsten U, Putnam N, et al. 2012. Phytozome: A comparative platform for green plant
941    genomics. *Nucleic Acids Res.* 40:1178–1186.

942 Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,
943    Raychowdhury R, Zeng Q et al. 2013. Trinity: reconstructing a full-length transcriptome
944    without a genome from RNA-Seq data. *Nat. Biotechnol.* 29:644–652.

945 Grant PR, Grant BR. 2014. Synergism of Natural Selection and Introgression in the Origin of
946    a New Species. *Am. Nat.* 183:671–681.

947 Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, Patterson N, Li H, Zhai W,
948    Fritz MH, et al. 2010. A Draft Sequence of the Neandertal Genome. *Science* 328:710–
949    722.

950 Gutjahr C, Gobbato E, Choi J, Riemann M, Johnston MG, Summers W, Carbonnel S,
951    Mansfield C, Yang S, Nadal M, et al. 2015. Rice perception of symbiotic arbuscular
952    mycorrhizal fungi requires the karrikin receptor complex. *Science* 350: 1521-1524.

953 Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB,
954    Eccles D, Li B, Lieber M, et al. 2013. De novo transcript sequence reconstruction from
955    RNA-seq using the Trinity platform for reference generation and analysis. *Nat. Protoc.*
956    8:1494–1512.

957 Hadid Y, Pavlicek T, Beiles A, Ianovici R, Raz S, Nevo E. 2014. Sympatric incipient
958    speciation of spiny mice *Acomys* at "Evolution Canyon," Israel. *Proc. Natl. Acad. Sci.*
959    111:1043–1048.

960 Hadid Y, Tzur S, Pavlicek T, Sumbera R, Skliba J, Lovy M, Fragman-Sapir O, Beiles A, Arieli
961    R, Raz S, et al. 2013. Possible incipient sympatric ecological speciation in blind mole
962    rats (*Spalax*). *Proc. Natl. Acad. Sci*. 110:2587–2592.

963    Hipperson H, Dunning LT, Baker WJ, Butlin RK, Hutton I, Papadopulos AST, Smadja CM,

964         Wilson TC, Devaux C, Savolainen V. 2016. Ecological speciation in sympatric palms: 2.

965         Pre- and post-zygotic isolation. *J. Evol. Biol.* 29:2143–2156.

966    Jombart, T, Ahmed I. 2011. Adegenet 1.3-1: new tools for the analysis of genome-wide SNP

967         data. *Bioinformatics* 27:3070–3071.

968    Kannan S, Hui J, Mazooji K. 2016. Shannon: An Information-Optimal de Novo RNA-Seq

969         Assembler. 2016. *bioRxiv* https://www.biorxiv.org/content/10.1101/039230v1

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

971         sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* 30:3059–3066.

972    Kautt AF, Machado-Schiaffino G, Torres-Dowdall J, Meyer A. 2016. Incipient sympatric

973         speciation in Midas cichlid fish from the youngest and one of the smallest crater lakes in

974         Nicaragua due to differential use of the benthic and limnetic habitats? *Ecol. Evol.*

975         6:5342–5357.

976    Kirkpatrick M. 2001. Reinforcement during ecological speciation. *Proc. R. Soc. B Biol. Sci.*

977         268:1259–1263.

978    Kondrashov AS, Kondrashov FA. 1999. Interactions among quantitative traits in the course

979         of sympatric speciation. *Nature* 400:351–354.

980    Levin DA. 2009. Flowering-time plasticity facilitates niche shifts in adjacent populations. *New*

981         *Phytol.* 183:661–666.

982    Levis NA, Isdaner AJ, Pfennig DW. 2018. Morphological novelty emerges from pre-existing

983         phenotypic plasticity. *Nat. Ecol. Evol.* 2:1289–1297.

984    Li H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping

985         and population genetical parameter estimation from sequencing data. *Bioinformatics*

986         27:2987–2993.

987    Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler

988         transform. *Bioinformatics* 25:1754–1760.

989    Linklater M, Hamylton S, Brooke B, Nichol S, Jordan A, Woodroffe C. 2018. Development of

990         a seamless, high-resolution bathymetric model to compare Reef morphology around the

991         subtropical island shelves of Lord Howe Island and Balls Pyramid, southwest Pacific

992         Ocean. *Geosciences* 8:11.

993    Liu J, Li G, Chang Z, Yu T, Liu B, Mcmullen R. 2016. BinPacker: Packing-Based De Novo

994         Transcriptome Assembly from RNA-seq Data. *PLoS Comput. Biol.* 12:1004772.

995    Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for

996         RNA-seq data with DESeq2. *Genome Biol.* 15:550.

997    Macmanes MD. 2014. On the optimal trimming of high-throughput mRNA sequence data.

998         *Front. Genet.* 5:13.

999    Malinsky M, Challis RJ, Tyers AM, Schiffels S, Terai Y, Ngatunga BP, Miska EA, Durbin R,
1000         Genner MJ, Turner GF. 2015. Genomic islands of speciation separate Cichlid
1001         ectomorphs in an East African crater lake. *Science* 350:1493–1498.
1002   Martin CH, Cutler JS, Friel JP, Touokong CD, Coop G, Wainwright PC. 2015. Complex
1003         histories of repeated gene flow in Cameroon crater lake cichlids cast doubt on one of
1004         the clearest examples of sympatric speciation. *Evolution* 69:1406–1422.
1005   Martin SH, Davey JW, Jiggins CD. 2015. Evaluating the use of ABBA-BABA statistics to
1006         locate introgressed loci. *Mol. Biol. Evol.* 32:244–257.
1007   McDougall I, Embleton BJJ, Stone DBI. 1981. Origin and evolution of Lord Howe Island,
1008         Southwest Pacific Ocean. *J. Geol. Soc. Aust.* 28:155–176.
1009   Meier JI, Marques DA, Mwaiko S, Wagner CE, Excoffier L, Seehausen O. 2017. Ancient
1010         hybridization fuels rapid cichlid fish adaptive radiations. *Nat. Commun.* 8:1–11.
1011   Miller KG, Wright JD, Browning J V, Kulpecz A, Kominz M, Naish TR, Cramer BS, Rosenthal
1012         Y, Peltier WR, Sosdian S. 2012. High tide of the warm Pliocene: Implications of global
1013         sea level for Antarctic deglaciation. *Geology* 40:407–410.
1014   Mirarab S, Warnow T. 2015. ASTRAL-II: Coalescent-based species tree estimation with
1015         many hundreds of taxa and thousands of genes. *Bioinformatics* 31:44–52.
1016   Nei M. 1987. Molecular Evolutionary Genetics. New York: Columbia University Press
1017   Nosil P, Feder JL, Flaxman SM, Gompert Z. 2017. Tipping points in the dynamics of
1018         speciation. *Nat. Ecol. Evol.* 1:1–8.
1019   Osborne OG, Batstone TE, Hiscock SJ, Filatov DA. 2013. Rapid speciation with gene flow
1020         following the formation of Mt. Etna. *Genome Biol. Evol.* 5:1704–1715.
1021   Osborne OG, De-Kayne R, Bidartondo MI, Hutton I, Baker WJ, Turnbull CGN, Savolainen V.
1022         2018. Arbuscular mycorrhizal fungi promote coexistence and niche divergence of
1023         sympatric palm species on a remote oceanic island. *New Phytol.* 217:1254–1266.
1024   Palmer AR. 2014. Symmetry Breaking and the Evolution of Development. *Science* 306:828–
1025         833.
1026   Papadopulos AST, Baker WJ, Crayn D, Butlin RK, Kynast RG, Hutton I, Savolainen V. 2011.
1027         Speciation with gene flow on Lord Howe Island. *Proc. Natl. Acad. Sci.* 108:13188–
1028         13193.
1029   Papadopulos AST, Price Z, Devaux C, Hipperson H, Smadja CM, Hutton I, Baker WJ, Butlin
1030         RK, Savolainen V. 2013. A comparative analysis of the mechanisms underlying
1031         speciation on Lord Howe Island. *J. Evol. Biol.* 26:733–745.
1032   Papadopulos AST, Igea J, Smith TP, Hutton I, Baker WJ, Butlin RK, Savolainen V. 2019a.
1033         Ecological speciation in sympatric palms: 4. Demographic analyses support speciation
1034         of *Howea* in the face of high gene flow. *Evolution* in press

1035 Papadopulos AST, Igea J, Dunning LT, Osborne OG, Quan X, Pellicer J, Turnbull C, Hutton
1036     I, Baker WJ, Butlin RK, Savolainen V. 2019b. Ecological speciation in sympatric palms:
1037     3. Genetic map reveals genomic islands underlying species divergence in *Howea*.
1038     *Evolution* in pressParadis E, Claude J, Strimmer K. 2004. APE: Analyses of
1039     phylogenetics and evolution in R language. *Bioinformatics* 20:289–290.
1040 Pardo-Diaz C, Baxter SW, Joron M, McMillan WO, Jiggins CD, Merot C, Salazar C,
1041     Figueiredo-Ready W. 2012. Adaptive Introgression across Species Boundaries in
1042     Heliconius Butterflies. *PLoS Genet.* 8:1002752.
1043 Pease JB, Hahn MW. 2015. Detection and Polarization of Introgression in a Five-Taxon
1044     Phylogeny. *Syst. Biol.* 64:651–662.
1045 Peng Y, Leung HCM, Yiu SM, Lv MJ, Zhu XG, Chin FYL. 2013. IDBA-tran: A more robust de
1046     novo de Bruijn graph assembler for transcriptomes with uneven expression levels.
1047     *Bioinformatics* 29:326–334.
1048 Peverill KI, Sparrow LA, Reuter DJ. 1999. Soil Analysis: An Interpretation Manual. Clayton,
1049     Australia: Csiro Publishing
1050 Pie, MR. The Influence of Phylogenetic Uncertainty on the Detection of Positive Darwinian
1051     Selection. 2006. *Mol. Biol. Evol.* 23, 2274-2278.
1052 Pfeifer B, Wittelsbürger U, Ramos-Onsins SE, Lercher MJ. 2014. PopGenome: An efficient
1053     swiss army knife for population genomic analyses in R. *Mol. Biol. Evol.* 31:1929–1936.
1054 Pfennig DW, Wund MA, Snell-Rood EC, Cruickshank T, Schlichting CD, Moczek AP. 2010.
1055     Phenotypic plasticity's impacts on diversification and speciation. *Trends Ecol. Evol.*
1056     25:459–467.
1057 Poelstra JW, Richards EJ, Martin CH. 2018. Speciation in sympatry with ongoing secondary
1058     gene flow and a potential olfactory trigger in a radiation of Cameroon cichlids. *Mol. Ecol.*
1059     27:4270-4288
1060 Przedpelska-Wasowicz E, Wasowicz P. 2013. Does zinc concentration in the substrate
1061     influence the onset of flowering in *Arabidopsis arenosa* (Brassicaceae)? *Plant Growth*
1062     *Regul.* 69:87–97.
1063 R Core Development Team. 2008. R: A language and environment for statistical computing.
1064     Vienna, Austria: R Foundation for Statistical Computing
1065 Riesch R, Muschick M, Lindtke D, Villoutreix R, Comeault AA, Farkas TE, Lucek K, Hellen E,
1066     Soria-Carrasco V, Dennis SR, et al. 2017. Transitions between phases of genomic
1067     differentiation during stick-insect speciation. *Nat. Ecol. Evol.* 1:0082.
1068 Robertson G, Schein J, Chiu R, Corbett R, Field M, Jackman SD, Mungall K, Lee S, Okada
1069     HM, Qian JQ, et al. 2010. De novo assembly and analysis of RNA-seq data. *Nat.*
1070     *Methods* 7:909–912.

1071 Rohlfs R V., Nielsen R. 2015. Phylogenetic ANOVA: The expression variance and evolution
1072     model for quantitative trait evolution. *Syst. Biol.* 64:695–708.

1073 Ryan PG, Bloomer P, Moloney CL, Grant TJ, Delport W. 2007. Ecological Speciation in
1074     South Atlantic Finches. *Science* 315:1420–1422.

1075 Sanderson MJ. 2003. r8s: inferring absolute rates of molecular evolution and divergence
1076     times in the absence of a molecular clock. *Bioinformatics* 19:301–302.

1077 Savolainen V, Anstett M-C, Lexer C, Hutton I, Clarkson JJ, Norup M V, Powell MP,
1078     Springate D, Salamin N, Baker WJ. 2006. Sympatric speciation in palms on an oceanic
1079     island. *Nature* 441:210–213.

1080 Schliewen UK, Tautz D, Pääbo S. 1994. Sympatric speciation suggested by monophyly of
1081     crater lake cichlids. *Nature* 368:629–632.

1082 Schulz MH, Zerbino DR, Vingron M, Birney E. 2012. Oases: Robust de novo RNA-seq
1083     assembly across the dynamic range of expression levels. *Bioinformatics* 28:1086–1092.

1084 Seehausen O. 2004. Hybridization and adaptive radiation. *Trends Ecol. Evol.* 19:198–207.

1085 Seehausen O. 2015. Process and pattern in cichlid radiations - inferences for understanding
1086     unusually high rates of evolutionary diversification. *New Phytol.* 207:304–312.

1087 Shimodaira H, Hasegawa M. 1999. Multiple comparisons of log-likelihoods with applications
1088     to phylogenetic inference. *Mol. Biol. Evol.* 16:1114-1116.

1089 Simao FA, Waterhouse RM, Ioannidis P, Kriventseva E V., Zdobnov EM. 2015. BUSCO:
1090     Assessing genome assembly and annotation completeness with single-copy orthologs.
1091     *Bioinformatics* 31:3210–3212.

1092 Song L, Florea L. 2015. Rcorrector: efficient and accurate error correction for Illumina RNA-
1093     seq reads. *Gigascience* 4:48.

1094 Sorenson MD, Sefc KM, Payne RB. 2003. Speciation by host switch in brood parasitic
1095     indigobirds. *Nature* 424:928–931.

1096 Stamatakis A. 2014. RAxML version 8: A tool for phylogenetic analysis and post-analysis of
1097     large phylogenies. *Bioinformatics* 30:1312–1313.

1098 Stankowski S, Streisfeld MA. 2015. Introgressive hybridization facilitates adaptive
1099     divergence in a recent radiation of monkeyflowers. *Proc. R. Soc. B Biol. Sci.* 282.

1100 Tajima F. 1989. Statistical method for testing the neutral mutation hypothesis by DNA
1101     polymorphism. *Genetics* 123:585–595.

1102 Wang B, Jin SH, Hu HQ, Sun YG, Wang YW, Han P, Hou BK. 2012. UGT87A2, an
1103     *Arabidopsis* glycosyltransferase, regulates flowering time via FLOWERING LOCUS C.
1104     *New Phytol.* 194:666–675.

1105 Wang WY, Xu J, Liu XJ, Yu Y, Ge Q. 2012. Cadmium induces early flowering in *Arabidopsis.*
1106     *Biol. Plant.* 56:117–120.

1107     Weir BS, Cockerham CC. 1984. Estimating F-Statistics for the Analysis of Population

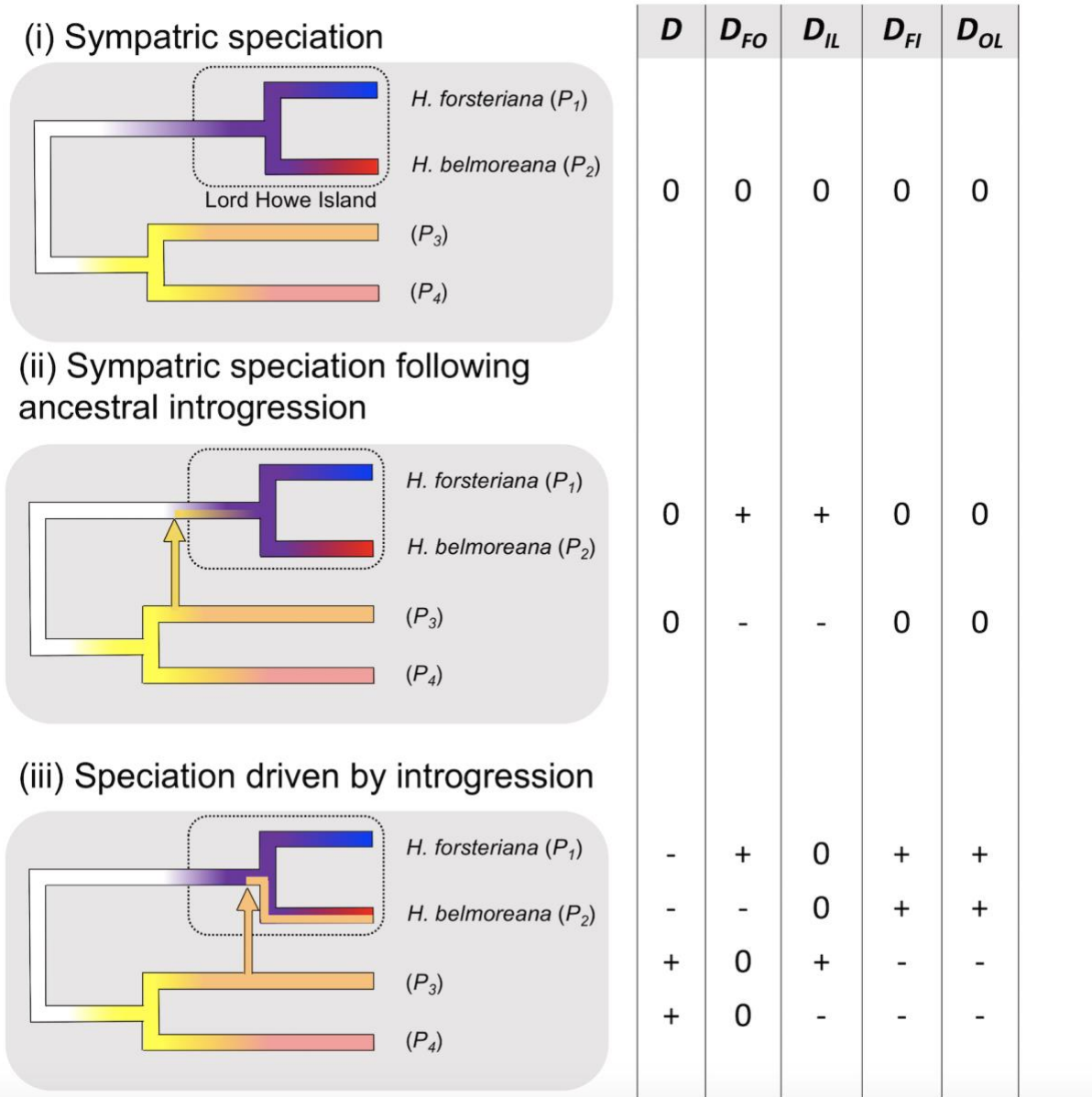1108         Structure. *Evolution* 38:1358–1370.

1109     Woodroffe CD, Kennedy DM, Brooke BP, Dickson ME. 2006. Geomorphological evolution of

1110         Lord Howe Island and carbonate production at the latitudinal limit to reef growth. *J*

1111         *Coast Res* 22:188–201.

1112     Xie Y, Wu G, Tang J, Luo R, Patterson J, Liu S, Huang W, He G, Gu S, Li S, et al. 2014.

1113         SOAPdenovo-Trans: De novo transcriptome assembly with short RNA-Seq reads.

1114         *Bioinformatics* 30:1660–1666.

1115     Yang Z. 2007. PAML 4: Phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.*

1116         24:1586–1591.

1117     Zhang C, Rabiee M, Sayyari E, Mirarab S. 2018. ASTRAL-III: Polynomial time species tree

1118         reconstruction from partially resolved gene trees. *BMC Bioinformatics* 19:15–30.

1119     Zhang J, Nielsen R, Yang Z. 2005. Evaluation of an improved branch-site likelihood method

1120         for detecting positive selection at the molecular level. *Mol. Biol. Evol.* 22:2472–2479.

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130



1131

Figure 1. Three possible scenarios of introgression during the evolution of *Howea*, with expected *D*-statistic results for each scenario: (i) ancestral *Howea* speciates allopatrically from the ancestor of *Linospadix* and *Laccospadix*, and later speciate sympatrically on LHI. In this case, all *D* statistics should not be significantly different from zero; (ii) sympatric speciation in *Howea* follows ancestral introgression between *Howea* and *Linospadix* or *Laccospadix* lineages. In this case, *D* should be zero, $D_{FO}$ and $D_{IL}$ should either both be positive (in the case of $P_3$ introgression) or both negative (in the case of $P_4$ introgression) whereas $D_{FI}$ and $D_{OL}$ should both be zero; (iii) *Howea* speciation is a direct result of introgression. In this case, a wide range of combinations of $D_{FOIL}$ statistics are possible (see Pease and Hahn, 2015, for details) but *D* will always be significantly positive or negative. The arrows on the phylogenies represent introgression events, colour change represents allele frequency changes over time
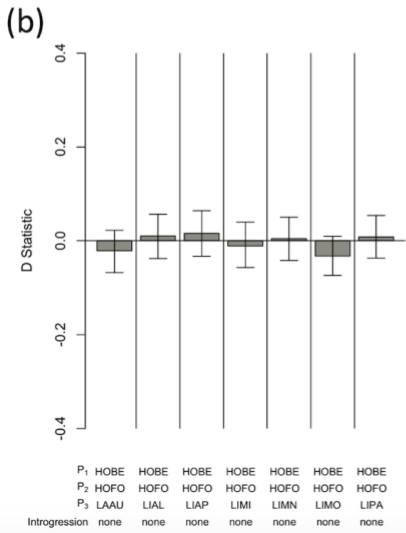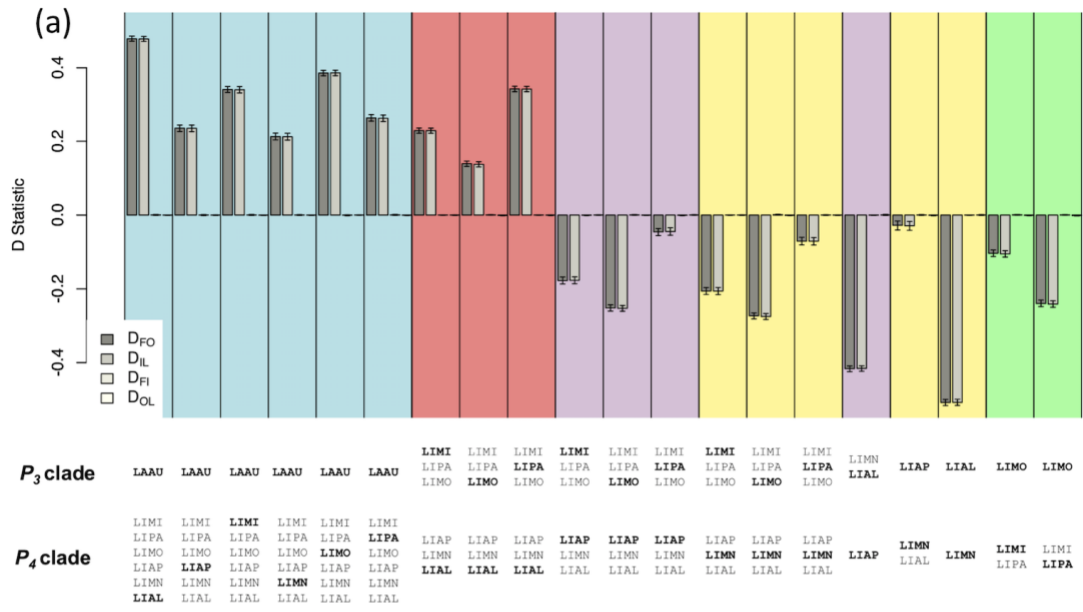
and the dotted boxes represents LHI. The table on the right shows the expected sign for Patterson's $D$ (Green et al. 2010) and $D_{FOIL}$ statistics ($D_{FO}$, $D_{IL}$, $D_{FI}$ and $D_{OL}$) for each scenario, with multiple possible combinations supporting some scenarios (Following Pease and Hahn 2015). While not represented in this figure, non-zero $D$-statistics can also result from ancestral population structure (see discussion).

Figure 2. Phylogeny and morphology of *Howea* and its closest relatives. (a) A dated species tree of *Howea*, *Linospadix* and *Laccospadix*. Node labels show the percentage of gene trees supporting the dominant topology followed by the other two possible (unrooted) topologies. Nodes are coloured by proportion of gene trees supporting the dominant topology and blue node bars show 95% confidence intervals of node ages. Illustrations of selected species are drawn to the right of the phylogenetic tree, with a 0.5 m scale bars to show approximate relative heights of the species. A DensiTree plot (b) shows the level of concordance between gene trees. Following filtering for low confidence nodes, each unique gene tree topology is transparently plotted such that gene tree discordance is apparent. Coloured bars denote species and are coloured according to the boxes beside the species names on panel (a). Photographs of the crowns of *H. belmoreana* (c) and *H. forsteriana* (d) show differences in leaf morphology.
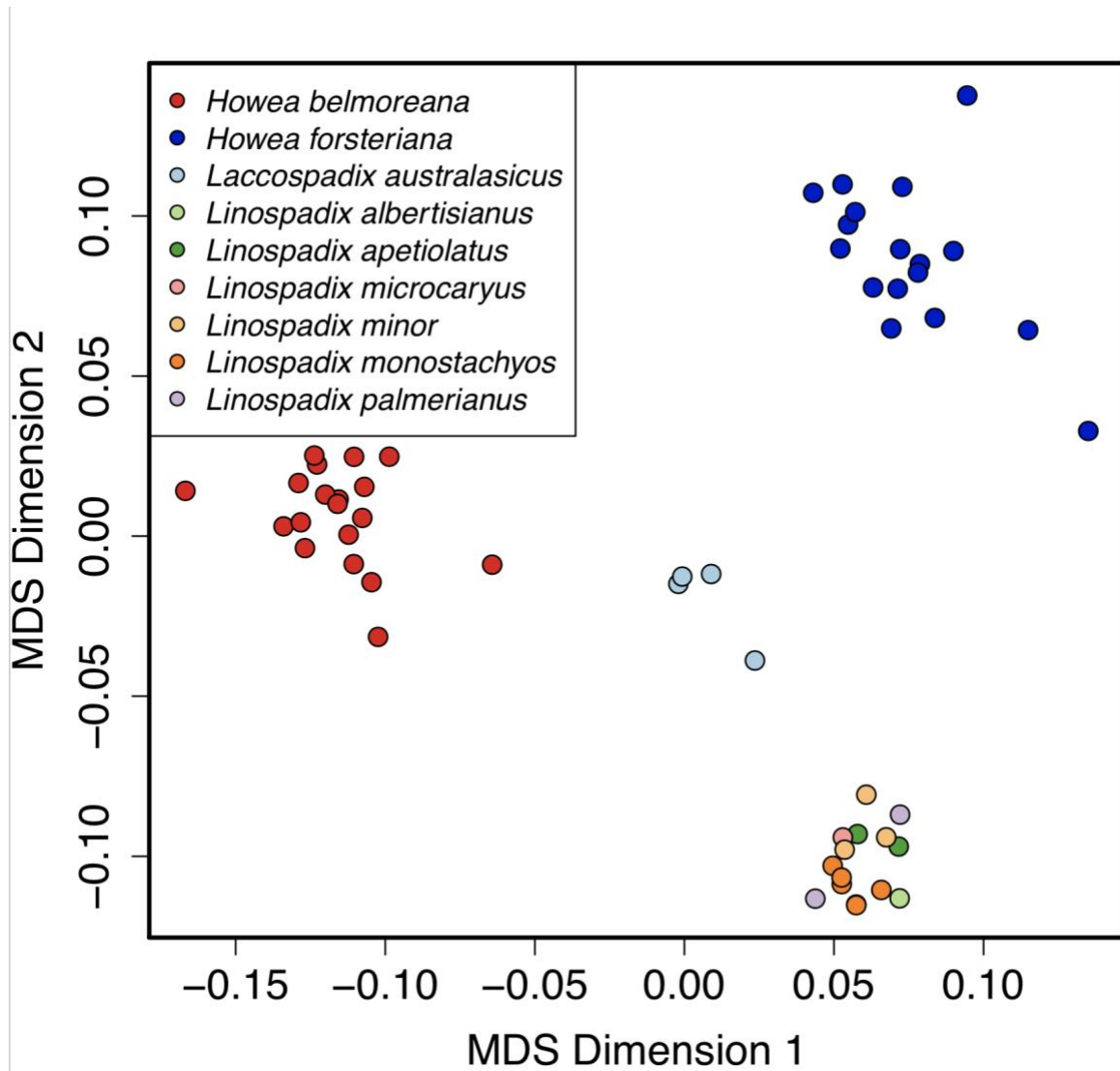
Figure 3. Introgression between taxa. Values and 95% confidence intervals are shown for all four $D_{FOIL}$ statistics (a), and Patterson's $D$ statistic (b) for each testable subtree. The identity of species in clades represented by $P_3$ and $P_4$ are shown below each test, with the actual species tested highlighted in bold. Tests unanimously support introgression between *Linospadix* or *Laccospadix* species and ancestral *Howea* (a), but not between extant species of *Howea* and any of the outgroups (b). Species abbreviations are as follows: HOBE: *Howea belmoreana*, HOFO: *H. forsteriana*, LAAU: *Laccospadix australasicus*, LIAL: *Linospadix albertisianus*, LIAP:*L. apetiolatus*, LIMI: *L. microcaryus*, LIMN: *L. minor*, LIMO: *L. monostachyos*, LIPA: *L. palmerianus*. One interpretation of the results which minimises the number of introgression events is shown in panel (c), with introgression events shown on the tree as coloured arrows corresponding to the tests that support them in panel (a). The horizontal positions of introgression arrows are arbitrary and do not reflect the
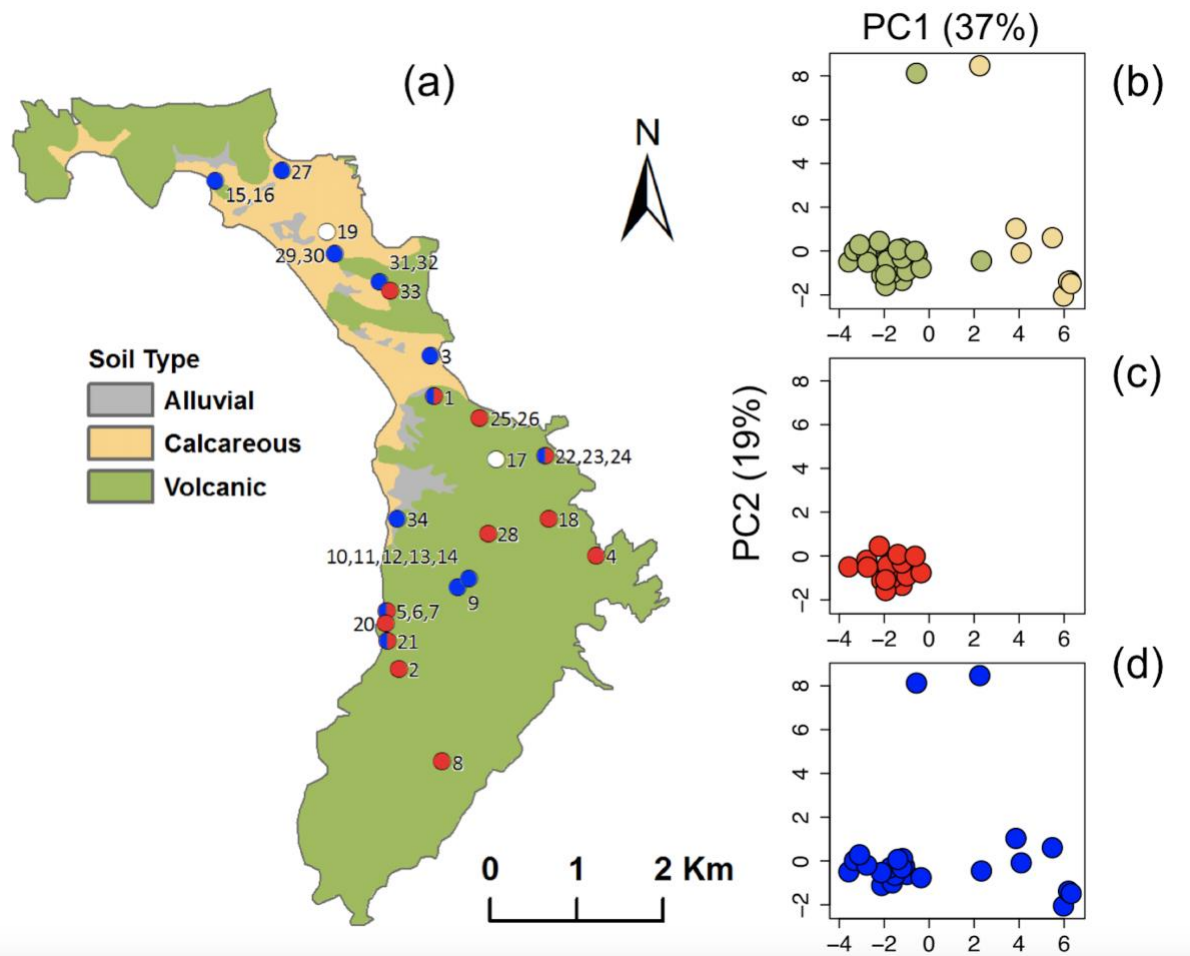
1176    timing of introgression. While we have interpreted the results in terms of introgression in

1177    panel (c), ancestral population structure can also lead to non-zero *D*-statistics.

1178

Figure 4. Genetic clustering amongst *Howea* and their outgroups revealed by
multidimensional scaling of all single nucleotide polymorphism data. The two *Howea* species
are equidistant from outgroup species, however *Laccospadix* is closer to *Howea* than its
sister genus *Linospadix* (see text).

Figure 5. Soil characteristics of the two *Howea* species habitats on LHI. (a) A map of the island showing broad soil classifications and sampling locations. Sampling sites with only *H. forsteriana* are shown in blue, with only *H. belmoreana* are shown in red and those with both species are shown half blue and half red. (b-d) The first two principle components (PCs) of a PCA of normalised soil metrics for water content, pH, concentrations of 20 acid extractable elements and four DTPA-extractable micronutrients. Plot (b) is coloured by soil type (volcanic: green, calcareous: yellow) and numbers on the plot correspond to those in (a). (c) shows all sites with *H. belmoreana* present (red) and (d) shows all sites with *H. forsteriana* present (blue).

1203

1204

1205

1206