# A blueprint of seed desiccation sensitivity in the genome of Castanospermum australe

Alexandre Marques, Maria-Cecilia D Costa, Jill M Farrant, Henk Hilhorst, Julia Buitink, Wilco Ligterink, Olivier Leprince, Sandra Pelletier, Toni Gabaldon, M. Eric Schranz, et al.

HAL Id: hal-02624811
https://hal.inrae.fr/hal-02624811

Submitted on 26 May 2020

1    **A blueprint of seed desiccation sensitivity in the genome of *Castanospermum australe***

2

3    **Alexandre Marques[1][†], Maria-Cecília D. Costa[2][†], Udisha Chathuri[2], Eef Jonkheer[3,4], Tao**

4    **Zhao[4], Elio Schijlen[5], Martijn Derks[3], Harm Nijveen[3], Marina Marcet-Houben[6,7], Irene**

5    **Julca[6,7], Julien Delahaie[8], M. Eric Schranz[4], Toni Gabaldon[6,7,9], Sandra Pelletier[8],**

6    **Olivier Leprince[8], Wilco Ligterink[1], Julia Buitink\*[8], Henk W.M. Hilhorst\*[1], Jill M.**

7    **Farrant\*[2]**

8    [1]Laboratory of Plant Physiology, Wageningen University and Research, Wageningen, The

9    Netherlands; [2]Department of Molecular and Cell Biology, University of Cape Town, Private

10    Bag, Rondebosch 7701, South Africa; [3]Bioinformatics Group, Wageningen University and

11    Research, Wageningen, The Netherlands; [4]Biosystematics Group, Wageningen University and

12    Research, Wageningen, The Netherlands; [5]Bioscience, Wageningen Plant Research

13    International, Wageningen, The Netherlands; [6]Centre for Genomic Regulation (CRG), The

14    Barcelona Institute of Science and Technology, Dr. Aiguader 88, Barcelona 08003, Spain;

15    [7]Universitat Pompeu Fabra (UPF), Barcelona 08003, Spain; [8]Institut de Recherche en

16    Horticulture et Semences, UMR1345, INRA, Agrocampus Ouest, Université d'Angers, SFR

17    4207 QASAV, 49071 Beaucouzé, France; [9]ICREA, Pg. Lluís Companys 23, Barcelona

18    08010, Spain.

19

20    [†]These authors contributed equally to this work

21    *Corresponding authors

22

23    Author for correspondence:

24    Henk W.M. Hilhorst

25    Tel: +31 317483646

26    Email: henk.hilhorst@wur.nl

27

28    **Summary**

29

30    •    Most angiosperms produce seeds that are desiccated on dispersal with the ability to

31    retain viability in storage facilities for prolonged periods. However, some species produce

32    desiccation sensitive seeds which rapidly lose viability in storage, precluding *ex situ*

33    conservation. Current consensus is that desiccation sensitive seeds either lack or do not

34    express mechanisms necessary for the acquisition of desiccation tolerance.

35    •    We sequenced the genome of *Castanospermum australe*, a legume species producing

36    desiccation sensitive seeds, and characterized its seed developmental physiology and -

37    transcriptomes.

38    •    *C. australe* has a low rate of evolution, likely due to its perennial life-cycle and long

39    generation times. The genome is syntenic with itself, with several orthologs of genes from

40    desiccation tolerant legume seeds, from gamma whole-genome duplication events being

41    retained. Changes in gene expression during development of *C. australe* seeds, as

42    compared to desiccation tolerant *Medicago truncatula* seeds, suggest they remain

43    metabolically active, prepared for immediate germination.

44    •    Our data indicates that the phenotype of *C. australe* seeds arose through few changes

45    in specific signalling pathways, precluding or bypassing activation of mechanisms

46    necessary for acquisition of desiccation tolerance. Such changes have been perpetuated as

47    the habitat in which dispersal occurs is favourable for prompt germination.

48

49    **Key words**: desiccation-sensitivity, evolution, orthodox, phylome, recalcitrant, seed

50    development, synteny

51

52    **Introduction**

53    Seeds of most gymnosperm and angiosperm species are shed in the desiccated state and can

54    be stored dry under sub-zero temperatures for prolonged periods of time, thus facilitating

55    plant germplasm conservation. Such desiccation tolerant (DT) seeds are termed 'orthodox'.

56    However, seeds of some species are desiccation-sensitive (DS), also referred to as

57    'recalcitrant', and cannot be successfully stored under typical conditions (Berjak &

58    Pammenter, 2013). DS-seeded species are mostly found in the humid tropics and may

59    represent up to 50% of the species present in tropical evergreen rain forests (Hamilton *et al.*,

60    2013). Such species occur in environments conducive to immediate seed germination and thus

61    selective pressure for desiccation tolerance has been relaxed or is absent. It has been

62    hypothesized that seed desiccation sensitivity is a derived trait that evolved independently in

63    non-related clades (Berjak & Pammenter, 2000). Genes responsible for seed desiccation

64    tolerance would have been lost, repressed and/or mutated in DS seeded species (Berjak &

65    Pammenter, 2008). However, this hypothesis remains to be tested at the genome level. Here

66    we applied an extensive phylogenetic comparison to obtain a genomic blueprint of desiccation

67    sensitivity in seeds.

68        Desiccation tolerance is acquired mid-way during the development of orthodox seeds

69    when seed filling is approximately half-way, corresponding to a steep drop in water content of

70    the seeds, concomitantly with a transient rise in abscisic acid (ABA) content (Bewley *et al.*,

71    2013). This acquisition comprises highly coordinated molecular events, including the

72    repression of photosynthesis and energy metabolism, and accumulation of protective

73    components, such as late embryogenesis abundant (LEA) proteins, anti-oxidants, and soluble

74    sugars (Leprince *et al.*, 2017). These events are tightly regulated by hormones such as ABA

75    and transcription factors (TFs) such as *ABSCISIC ACID INSENSITIVE 3* (*ABI3*), *FUSCA 3*

76    (*FUS3*) and *LEAFY COTYLEDON 1* (*LEC1*) (Leprince *et al.*, 2017). Conversely, in the

77    development of DS seeds, the acquisition of desiccation tolerance and accumulation of the

78    above-mentioned protectants appear to be suppressed and rather they directly progress

79    towards germination (Farrant *et al.*, 1993b; Francini *et al.*, 2006; Delahaie *et al.*, 2013).

80    However, the genetic makeup underlying the DS seed phenotype is unknown.

81        The legume family (Fabaceae) contains many agriculturally important species, all

82    producing DT seeds, such as soybean (*Glycine max*), the common bean (*Phaseolus vulgaris*),

83    lentils (*Lens culinaris*) and chickpeas (*Cicer arietinum*). Moreover, legume species, such as

84    soybean and *Medicago truncatula*, are important experimental models for molecular and

85    physiological studies on seed desiccation tolerance (Chatelain *et al.*, 2012; Delahaie *et al.*,

86    2013; Verdier *et al.*, 2013; Zinsmeister *et al.*, 2016). In addition, the genomes of 16 DT-

87    seeded legume species have been sequenced. Thus, a large amount of information is available,

88    allowing comparative analysis among them.

89        *Castanopermum australe* A. Cunn & C. Fraser ex Hook, also known as the Moreton

90    Bay Chestnut or Blackbean*,* is a leguminous tropical tree native to the east coast of Australia

91    and west Pacific islands. In contrast to most legume species, *C. australe* produces DS seeds. It

92    is the only known species in the genus that forms a separate and early branching clade within

3

93    the Papilionoideae subfamily (Cardoso *et al.*, 2012). The earliest-branching papilionoids fall

94    within an ADA clade, which includes the monophyletic tribes Angylocalyceae, Dipterygeae,

95    and Amburanae. *C. australe*'s important phylogenetic position in the basis of the ADA clade

96    (Angylocalyaceae) makes its genome ideal to study both trait evolution and the ancient

97    polyploid history of papilionoid legumes (Schranz *et al.*, 2012).

98         Here, we provide detailed genomic sequence information of *C. australe* combined

99    with time-resolved gene expression analysis of seed development of this species, including

100   the comparison with other species producing either DS or DT seeds. Such information is key

101   to understanding mechanisms of desiccation tolerance and, ultimately, to design strategies to

102   improve tolerance of extreme water loss in DS seeds for conservation purposes. We

103   investigated genomic changes associated with seed desiccation sensitivity, including gene

104   deletions, severe mutations and gene mis-expression, as well as their relationship with gene

105   expression patterns during seed development and maturation.

106

107   **Materials and Methods**

108   Plant material

109   A population of trees of *C. australe* growing in Pietermaritzburg (Kwazulu-Natal Province,

110   South Africa) was the source of plant material for this work. Seed development occurs over a

111   6-month period, during which pods were harvested weekly. Seeds were extracted and,

112   following histodifferentiation, were separated into component tissues (axis, cotyledons and

113   seed coat). The following was determined annually over 2 seasons. Whole seed mass and that

114   of component tissues (n = 60) and water content (n=10-20) was determined gravimetrically by

115   oven drying. The ability of intact seeds to germinate was tested by planting in vermiculite.

116   Once germinable, the amount of water loss tolerated by axes and cotyledons was determined

117   by flash drying (Berjak *et al.*, 1990). Axis survival was determined as the ability to produce

118   both shoots and roots when cultured in vitro on full strength MS medium. Survival of

119   cotyledons was assessed by tetrazolium staining followed by spectrometric analysis (Sershen

120   *et al.*, 2012). Critical water contents (calculated on a $gH_2O.g^{-1}$ dry mass) were calculated at

121   those stages at which 50% survival was observed.

122

123   Sugar and ABA content and determination

4

Comment citer ce document :
Marques, Costa, M.-C. D., Farrant, J. M., Hilhorst, Buitink, J., Ligterink, Leprince, O.,
Pelletier, S., Gabaldon, T., Schranz, M. E., Delahaie, J., Julca, I., Marcet-Houben, Nijveen, H.,
Derks, Schijlen, E., Zhao, Jonkheer, Chaturi (2019). A blueprint of seed desiccation sensitivity
in the genome of Castanospermum australe. BioRxiv, preprint. , DOI : 10.1101/665661

124    Sugars were extracted from frozen and lyophilised seeds and analysed by HPLC on a

125    Carbopac PA-1 column (Rosnoblet *et al.*, 2007) (Dionex Corp., Sunnyvale, CA, USA). Three

126    independent extractions and assays were performed on approx. 100 mg of tissue.

127    ABA was extracted and quantified as described by (Floková *et al.*, 2014).

128

129    Genome sequencing and assembly

130    Freeze-dried leaf material was used for DNA isolation as described by (Bernatzky &

131    Tanksley, 1986) with modifications. The genomic *C. australe* library consisted of 30x

132    coverage PacBio with a mean read length of 7.8 Kb. In addition, an Illumina paired end

133    library with reads of 100 bp and a 200-400 bp insert size was constructed and sequenced to a

134    64x coverage. Reads originating for contaminants were removed from all sequence data prior

135    to assembly. Organelle genomes were also removed from the main genome assembly.

136    Illumina reads were error-corrected using Lighter (Song *et al.*, 2014) and assembled using

137    SparseAssembler (Ye *et al.*, 2012). A hybrid assembly was produced with DBG2OLC (Ye *et

138    al.*, 2016) and the contigs were reordered and connected into scaffolds using SSPACE-

139    LongRead (Boetzer & Pirovano, 2014). The assembly was polished using Sparc (Ye *et al.*,

140    2012) and Pilon (Walker *et al.*, 2014). PBJelly2 (English *et al.*, 2012) was used for gap

141    closure and genome improvement. Alignments due to gene duplication and repeats were

142    filtered out using the delta-filter utility of the MUMmer package (Kurtz *et al.*, 2004). The

143    assembly was validated by mapping the available RNA and DNA libraries to the genome with

144    Bowtie2 (Langmead & Salzberg, 2012) and Blasr (Chaisson & Tesler, 2012). Assembly

145    statistics were calculated using QUAST (Gurevich *et al.*, 2013). Gene space completeness

146    was measured using BUSCO (Benchmarking Universal Single-Copy Orthologs (Simão *et al.*,

147    2015)). The MAKER2 annotation pipeline (Holt & Yandell, 2011) was applied for gene

148    prediction and repeat annotation.

149

150    Transcriptome analysis

151    *C. australe* seeds were harvested at weekly intervals after seed set and at shedding from trees

152    growing in Pietermaritzburg in 2009 and 2011. Six developmental stages were collected for

153    cotyledon tissues, and three stages for the embryonic axes. Prior to maturation, when embryos

154    were small, cotyledon tissue was used. For the green pod stage (Figure 2), axes and

155    cotyledons were separated. Transcriptome analysis was performed on newly designed

156    12x135K Nimblegen arrays (IRHS_Ca_102K_v1) for *C. australe*, based on the genome

5

157  assembly of (Delahaie *et al.*, 2013). RNA amplification, labelling and hybridization was

158  performed according to (Terrasson *et al.*, 2013). Four biological replicates were analysed per

159  developmental stage using the dye-swap method, and statistical analysis on the gene

160  expression data was performed according to (Verdier *et al.*, 2013). Data are deposited in the

161  NCBI Gene Expression Omnibus database (accession no. GSE109217, samples

162  GSM2935461-GSM2935508). A gene was considered differentially expressed if $P \leq 0.05$ in

163  at least one comparison (axis or cotyledon) after the application of linear modelling.

164  Over-representation analysis (ORA) was used to recover over-represented biological

165  processes using the app BiNGO (default settings) (Maere *et al.*, 2005) for Cytoscape.

166  Orthologs were defined as hits with lowest Expect value (E-value) observing a

167  threshold of $<10^{-10}$. Multiple hits were considered orthologs when the difference between

168  their E-values and the lowest hit's E-value was smaller than $10^{-10}$.

169  Members of the eight LEA protein families were identified uploading Hidden Markov

170  Models (HMM) for each family from the PFAM database (Finn *et al.*, 2011) to HMMER

171  3.1b2 package (Eddy, 2011). All proteins with significant hits (E-value $\leq 0.01$) were selected.

172

173  Phylome reconstruction

174  We reconstructed three phylomes, one for a species set closely related to *C. australe*

175  (phylome 110) and starting in *C. australe*, a second one based on a broader taxonomic focus

176  also starting in *C. australe* (111) and a third one also broad but starting in *Fragaria vesca*

177  (112). The phylome 112 was used to search for lost genes in *C. australe*. It starts in *F. vesca*

178  that is the closest outgroup DT seeded species with a high-quality genome sequence. Both

179  phylomes were reconstructed using the same approach (Huerta-Cepas & Gabaldon, 2011). All

180  data generated during the phylome reconstruction has been deposited in phylomeDB

181  (PMID:24275491) under the phylomeID codes 110 and 111. The trees, alignments and

182  orthology and paralogy predictions are accessible to browse or download at the PhylomeDB

183  database (Huerta-Cepas *et al.*, 2014). A set of 183 one-to-one orthologous proteins present

184  across the compared species was used to reconstruct a species phylogeny.

185

186  Whole genome duplication analysis

187  The *C. australe* genome was compared to genomic data from five Legumes: *Glycine max*,

188  *Lotus japonicus*, *Medicago truncatula*, *Phaseolus vulgaris*, *Trifolium pratense*. *G. max* and

6

189  *M. truncatula* are tetraploids and the rest are diploid species. The diploid *Fragaria vesca*

190  species, of the family Rosaceae, was selected as outgroup as it is one of the closest relatives

191  with a completed genome available outside of the Legume family. The assemblies of: *G. max*,

192  *L. japonicus*, and *M. truncatula* are on chromosome level which made it easier to identify

193  genome collinearity and duplication patterns.

194

195  Synteny analysis

196  A synteny network approach (Zhao *et al.*, 2017) was implemented to compare the synteny of

197  *C. australe* and other whole-genome sequenced legume species available at the Legume

198  Information System (https://legumeinfo.org). BLASTP was used for pairwise genome

199  comparison. MCScanX (Wang *et al.*, 2012) was used for synteny block detection. Infomap

200  algorithm (Rosvall & Bergstrom, 2008) implemented in R igraph package was used for

201  synteny network clustering. Clusters containing genes/nodes from more than 8 (out of the 10)

202  legume species but no *C. australe* node(s) were screened out for further investigation. The

203  maximum distance between two matches was 20 genes, a syntenic block consists of minimum

204  5 genes and no blocks were merged. Quota Align was enabled to determine the syntenic depth

205  (the number of times a genomic region is syntenic). For calculating the fractionation bias, the

206  window size was lowered from 100 to 25 considering the smaller contigs of the *C. australe*

207  genome. Synonymous (Ks) and non-synonymous (Kn) site mutations were calculated for each

208  syntenic gene pair. Mutation rates were used to determine if genes were duplicated by the

209  WGD event or not. The distribution of the Ks rate was used to set a different cut-off per

210  species.

211

212  dN/dS analysis

213  We used the genome of 5 legume species from phytozome (*M. truncatula, G. max, Glycine*

214  *soja, T. pratense* and *P. vulgaris*) and *C. australe*. The SynMap tool in the online CoGe portal

215  was used to find syntenic gene pairs within these species. We set the DAGChainer on a

216  maximum distance between two matches for 50 genes and minimum number of aligned pairs

217  on 3 genes. To establish sets of orthologous among the 5 species against *M. truncatula*, the

218  method of reciprocal best hits using Last was used.

219      Codeml in the Phylogenetic Analysis by Maximum Likelihood (PAML) package was

220  used to estimate the dN (the rate of non-synonymous substitutions), dS (the rate of

7

221 synonymous substitutions) and the ratio of dN/dS (Yang, 2007). Orthologs with dS>5, dN>2

222 or dN/dS>2 were filtered. For genes with multiple syntelogs we kept the pair with the lowest

223 dN/dS.

224

225 **Results**

226 *Castanospermum australe* seed development

227 Seed development in *C. australe* occurs over a period of 6 months, with reserve accumulation

228 and maximum embryo size completed approximately 3 months after flowering and coincident

229 with pods becoming yellow (Figure 1 and Table 1). Unlike the embryo, the seed coat declined

230 in mass until just prior to the yellow-green pod stage, with no further decline once embryos

231 reached full size. Water content declined in all tissues once reserve accumulation was

232 complete at the yellow pod stage. This loss stabilized in all tissues but the seed coat, which

233 continued to lose water.

234 While seeds had the capacity to germinate prior to reaching full embryonic size,

235 germination rate was slow. Full germination capacity, typically of newly shed seeds was

236 achieved at the yellow green pod stage (Table 1). The lethal water content of axes below

237 which 50% viability was lost after drying declined from 0.45 $gH_2O/g^{-1}$ dry mass in those

238 extracted from green pods, to 0.23 $gH_2O/g^{-1}$ dry mass in axes from yellow pods, after which

239 there was no significant change. Cotyledons were more sensitive to dehydration, with 50%

240 loss of viability occurring below 0.82 $gH_2O/g^{-1}$ dry mass in those from green pods, this

241 declining slightly with development to 0.70 $gH_2O/g^{-1}$ dry mass in cotyledons from brown

242 pods (Table 1).

243 ABA regulates many aspects of plant growth and development including embryo

244 maturation, seed dormancy, germination, cell division and elongation. ABA content of *C.*

245 *australe* embryos was high during early developmental stages but declined considerably in

246 both axes and cotyledons in the transition from the yellow-green to the yellow pod stage

247 (Table 1). ABA content increased considerably in the seed coat in the transition from the

248 yellow to the brown pod stage, especially in the point of attachment (tissue that attaches the

249 embryo to the pod).

250

8

251     Seed maturation drying

252     DT seeds become tolerant of drying midway during seed development, concomitant with

253     reserve accumulation (Chatelain *et al.*, 2012). From this stage onwards, there is progressive

254     loss of water characterizing a process termed 'maturation' drying, that occurs after reserve

255     accumulation is complete and full seed size is attained. There was no maturation drying

256     typical of DT seeds after the yellow pod stage in *C. australe* (Table 1).

257         Previous work has identified transcripts that accumulate during the acquisition of

258     desiccation tolerance in *M. truncatula* seeds (Terrasson *et al.*, 2013; Righetti *et al.*, 2015).

259     Homologs of 121 of these genes failed to accumulate transcripts to a similar extent in *C.*

260     *australe* cotyledons in comparable seed developmental stages (Table S1). These transcripts

261     are related to sugar metabolism, photosynthesis, seed development, protection against abiotic

262     stress and modulation of plant stress responses. Examples include *ABI3*, *ABI5*, chaperone

263     proteins, heat shock factor proteins, putative LEAs, transparent testa protein, oleosins, *1-*

264     *CYSTEINE PEROXIREDOXIN*, and α-galactosidases.

265         We identified 269 transcripts with decreasing abundance in *C. australe* during final

266     maturation and increasing abundance in *M. truncatula* (Table S2). A certain number of these

267     transcripts are possibly involved in longevity (life span in the dried state) (Verdier *et al.*,

268     2013; Righetti *et al.*, 2015). Some of these genes are related to metabolic and catabolic

269     processes, such as *lipid metabolic process, cellular lipid metabolic process* and *nitrogen*

270     *compound metabolic process* (Table S3), reflecting the type of reserves accumulated. Mature

271     *C. australe* seeds predominantly accumulate starch (85% dry mass, Table1). Lipids and

272     proteins constitute only 3-6% and 2.8% dry mass, respectively. Mature seeds of *M. truncatula*

273     predominantly accumulate protein (30-40% dry mass), but also contain storage lipids (7-9%

274     dry mass) and a small amount of starch (< 1% dry mass) (Djemel *et al.*, 2005).

275         We also identified 296 transcripts with decreasing abundance in *M. truncatula* and

276     increasing abundance in *C. australe*. They are related to *root development*, *developmental*

277     *process* and *regulation of localization* (Table S2), which are likely associated with

278     germination related processes.

279         LEA proteins have been related to survival in the dry state (Chatelain *et al.*, 2012) and

280     to responses to environmental stresses, including desiccation (Tunnacliffe & Wise, 2007). In

281     *C. australe*, several LEA proteins failed to accumulate in cotyledons (Delahaie *et al.*, 2013).

282     A genome-wide search for LEAs in the genome of *C. australe* identified 94 LEA motif-

283     containing proteins, a number similar to what has been described for DT-seeded species, for

9

Version preprint

284    example, 88 in *Arabidopsis thaliana* and 99 in *Sorghum* bicolor, as well as for vegetative DT

285    in the resurrection plants *Oropetium thomaeum* and *Xerophyta viscosa* with 94 and 126

286    LEAs, respectively (VanBuren *et al.*, 2017; Costa *et al.*, 2017). A hierarchical clustering of

287    expression of the LEAs in developing seeds of *C. australe* and *M. truncatula* separated the

288    transcripts in two major clusters (Figure S1). LEAs in the first cluster belong to different

289    families and transcript abundance increased considerably in *M. truncatula* seeds towards mid

290    maturation. However, in *C. australe* most of these increased only slightly in early stages of

291    development but declined during the later stages. LEAs in the second cluster belong to the

292    LEA_2 family and do not undergo major changes in transcript abundance in either species.

293    This family encodes 'atypical' LEA proteins because of their more hydrophobic character

294    compared to other LEA families (Hundertmark & Hincha, 2008). Functional studies on

295    LEA_2 proteins suggest that they do not act in the protection of membranes in tissues

296    undergoing dehydration, although some proteins of this family were shown to have enzyme

297    protective properties under both freezing and drying conditions (Dang *et al.*, 2014).

298         Changes in soluble sugar content and composition have been described as a

299    characteristic of late maturation in DT seeds and correlate with the acquisition of longevity

300    and preparation for the dry state (Wang *et al.*, 2013; Leprince *et al.*, 2017). Whereas in all DT

301    legume seeds, raffinose family oligosaccharides (RFO) are the predominant sugars that

302    increase during late maturation, *C. australe* seeds were composed of 7-10% of soluble sugars,

303    with sucrose being the most abundant sugar detected and only minute amounts (0.7% of total

304    soluble sugars at the brown pod stage) of RFOs accumulating. Low ratios of sucrose:RFO

305    accumulation have been proposed to be a signature of desiccation tolerance and potential

306    indicators of seed storage categories (Steadman *et al.*, 1996; Farrant *et al.*, 2012) and the

307    opposite of this, as depicted in *C. australe* could be one of the reasons for desiccation

308    sensitivity in this species.

309    

310    Genome sequencing and assembly

311         In an effort to obtain a genomic blueprint of DS seeds we sequenced the genome of *C.*

312    *australe*. This species has a key position in the legume family phylogenetic tree, in the basis

313    of the ADA clade, which favours the study of trait evolution as well as the ancient polyploid

314    history of papilionoid legumes.

315         We produced an assembly with a total length of 382 Mb and an N50 of 832.6 Kb that

316    covers 96.7% of the predicted genome size. The assembly consisted of 1,210 contigs and

317    1,027 scaffolds (Table 2). The GC content was 32.9%. Genome annotation identified 29,124

318    protein-coding genes of which 98.1% show high sequence similarity to proteins in TrEMBL

319    and 84.4% in Swiss-Prot databases. An estimation of genome completeness indicated that

320    96.4% of the BUSCO (Benchmarking Universal Single-Copy Orthologs) genes were present.

321    Transposable elements covered 15.5% of the total genome. Repeat elements comprised 119

322    Kb of SINE (Short Interspersed Nuclear Element) retrotransposons, 383 Kb of LINE (Long

323    Interspersed Elements) retrotransposons, 13 Mb DNA transposons and 42 Mb of annotated

324    LTR (long terminal repeat-retrotransposons) sequences (Table 2).

325

326    Genomic alterations linked to seed desiccation sensitivity

327    To investigate the evolution of *C. australe* and legume diversification, a phylome was

328    constructed. A phylome constitutes the collection of all gene phylogenies in a genome. It is a

329    valuable source of information to establish evolutionary relationships among organisms and

330    their genes (Huerta-Cepas *et al.*, 2014). Our phylome contained the evolutionary histories of

331    all *C. australe* protein coding genes and their homologues in 20 publicly available sequenced

332    plant species (Figure 2). These species represent a broad phylogenetic distribution and include

333    the diverse seed storage phenotypes, i.e. DT, DS and intermediate seeds. Intermediate seeds

334    are typically tolerant of relatively extreme water loss, to 0.1-0.14 g $H_2O/g^{-1}$ dry mass but no

335    lower than this(Marques *et al.*, 2018), and have poor survival under conventional storage

336    conditions (Berjak & Pammenter, 2013). Such seeds thus show a storage phenotype in-

337    between orthodox and recalcitrant seeds.

338            The phylome analysis indicated that very few protein-coding genes ($\leq$ 1%) were

339    present in DS species only, and no protein-coding genes were retained in all DT species but

340    lost in all DS species (Figure 2). Several genes were lost in all DS species and retained in at

341    least half of the DT species (Table S4). 76 genes lost their ortholog but kept a paralog in *C.*

342    *australe* and 59 genes were lost in *C. australe* without detected paralogs, of which 11 were

343    shared with other DS-seeded species.

344            According to the phylome, 3,716 gene families were expanded exclusively in *C.*

345    *australe*, of which 180 were associated with transposons. Expansion size ranged from 2 to 32

346    genes and involved 30.5% of the predicted proteome. After removal of expansions associated

347    with transposable elements or viruses, the remaining expanded gene families were enriched

348    for GO (gene ontology) terms such as *defence response, flavonoid biosynthetic process*,

349    *terpene synthase activity* and *nutrient reservoir activity*. GO terms associated with *terpene*

11

350 *synthase activity, lyase activity* and *pectinesterase activity* were also enriched at the base of

351 the Papilionoideae subfamily. The phylome analysis also indicated that the duplication

352 frequency at the base of the Papilionoideae subfamily is lower (0.22) than that found at the

353 base of the Fabaceae family (0.73, Figure 3). Histograms of the synonymous rates and

354 average rates of syntenic blocks for six legume species showed distinctive peaks tracing back

355 to the shared papilionoid legume whole genome duplication (WGD), suggesting that the rate

356 of evolution is very diverse among closely related family members (Figure 4). The peak in *C.*

357 *australe* corresponds to a rate of synonymous (Ks) site mutations of 0.25 which is less than

358 half of the rate observed in *G. max* (0.6) and a third of that in *M. truncatula* (0.85). The

359 substitution rate in *C. australe* is so low that a second peak corresponding to the eudicot

360 hexaploidy (gamma WGD event) is still visible. The gamma event was also detected in the

361 histograms of block averages in *G. max* and *P. vulgaris*. The low Ks rate in *C. australe* is

362 likely due to this species being the only perennial in this comparative study and the one with

363 the longest generation times.

364      The evolution of a trait is shaped by the selective pressures to which it is subject.

365 Some selective pressures act to increase the benefits accumulated while others act to reduce

366 the costs incurred, affecting the cost/benefit ratio. Different selective pressures can be

367 estimated by the ratio of the number of nonsynonymous substitutions per non-synonymous

368 site (dN) in a specific period to the number of synonymous substitutions per synonymous site

369 (dS) in the same period (Mugal *et al.*, 2014).

370      Genome wide analysis of protein coding genes of *C. australe* in comparison with other

371 legume genomes enabled the identification of genes with 2-fold higher dN/dS in *C. australe*

372 (Table S5). Among these were genes associated with hormonal signalling, such as

373 *ACTIVATION-TAGGED BRI1* (*BRASSINOSTEROID-INSENSITIVE1*)-*SUPPRESSOR1*

374 (*ATBS1*), Ethylene Insensitive 3 family protein, *BRI1-ASSOCIATED RECEPTOR KINASE*

375 (*BAK1*) and *SALT TOLERANCE HOMOLOG2* (*STH2*). Other examples are: *SEEDSTICK*

376 (*STK*), *TRANSPARENT TESTA5* (*TT5*), *ENDO-BETA-MANNANASE7* (*MAN7), FLOWER*

377 *FLAVONOID TRANSPORTER* (*FFT*) and *ROTUNDIFOLIA3* (*ROT3*).

378      The evolution of a trait can also be investigated by analysing the degree to which

379 genes remain on corresponding chromosome (synteny) and in corresponding orders over time.

380 We investigated whether the loss of synteny in *C. australe* genes could be related to the loss

381 of seed desiccation tolerance. There are 169 genes re-arranged in the *C. australe* genome that

382 have syntenic orthologs in 50 angiosperm species. Most noteworthy among these were the

383 genes *BRASSINOSTEROID INSENSITIVE 3* and *5* (*BIN3* and *BIN5*), which participate in

12

384    brassinosteroid (BR) signalling and are associated with seed size (Yin *et al.*, 2002).

385    Furthermore, genes such as *MINISEED3* (*MINI3*) and *HAIKU1* (*IKU1*), regulators of seed

386    size via the BR pathway in *A. thaliana* (Luo *et al.*, 2005), also lost synteny in *C. austral*e,

387    which could contribute to the large seed size in this species.

388        The synteny between the genome of *C. australe* and other legume species was

389    evaluated by aligning the genome of *C. australe* against itself and against the genomes of *G.*

390    *max* and *M. truncatula* (Figure S2). Most of the *C. australe* genome is syntenic with itself and

391    mostly duplicated after the WGD event. While the duplicates are associated with the most

392    recent WGD event, many paralogs derived from the gamma event were also detected. The

393    alignment of *C. australe* against *G. max* indicated a high amount of syntenic orthologs and

394    paralogs, whereas the alignment of *C. australe* against *M. truncatula* indicated that although

395    many syntenic orthologs have been conserved, most of the WGD-derived paralogs were lost.

396    Moreover, most of the duplicated regions retained by *M. truncatula* were also duplicated and

397    retained in *C. australe*.

398

## Discussion

400    In orthodox seeds, survival in the dry state is a result of a series of molecular and cellular

401    processes that occur during the late stages of seed development. These processes result in the

402    acquisition of desiccation tolerance and longevity in the dry state. Although we have a

403    detailed understanding of these associated processes in orthodox seeds, limited information is

404    available regarding the development of seeds that do not fully activate them, such as *C.*

405    *australe*. Our study provides detailed information about *C. australe* seed development and

406    desiccation sensitivity.

407        In *C. australe* axes and cotyledons, some water loss occurss during reserve

408    accumulation but this stops once the embryo reaches its full size, stabilizing at 2.4 and 1.6

409    $gH_2O$ $g^{-1}$ dry mass, respectively (Table 1). This value is substantially higher than the 0.1

410    $gH_2O$ $g^{-1}$ dry mass reached by desiccation-tolerant seeds, such as *M. truncatula*.

411        The pattern of sugar accumulation also differs markedly between these species. The

412    percentage of soluble sugars in seeds of *C. australe* (7-10%) is comparable with the average

413    percentage for legume species (8-10% (Djemel *et al.*, 2005)). However, while RFOs are the

414    main sugars in *M. truncatula*, comprising 90% of the total soluble sugar content, in *C.*

415    *australe* only minute amounts of RFOs, mainly stachyose, could be detected (0.7% of total

416    soluble sugars at the brown pod stage). Stachyose and raffinose contents were highest in seeds

13

417    from the green pod stage and decreased with further progress of maturation (Table 1). A

418    similar finding has been reported for the non-viviparous highly DS seeds of *Avicennia marina*

419    (Farrant *et al.*, 1992). In parallel, sucrose and glucose content increased. The reduction of

420    stachyose content during further maturation suggests hydrolysis, normally occurring in

421    germinating DT seeds (Rosnoblet *et al.*, 2007). A comparative analysis of transcripts linked to

422    RFO metabolism between *C. australe* and *M. truncatula* identified transcripts of genes related

423    to the synthesis of sucrose from fructose-6 phosphate that remained high in *C. australe*

424    whereas they decreased in *M. truncatula*. Conversely, the transcripts of several genes related

425    to the synthesis of galactinol or raffinose and stachyose accumulated during development of

426    *M. truncatula* seeds while their abundance remained low in *C. australe* (Supplementary

427    Figure S1B). This set of genes might explain the lack of RFO accumulation in *C. australe*

428    seeds. Whereas the specific roles of RFOs in protection compared to the nonreducing sucrose

429    remain unconfirmed, the DS *A. thaliana abi3* mutants as well as Mt-*abi5* are also impaired in

430    the accumulation of RFOs (Zinsmeister *et al.*, 2016).

431       At the transcriptome level, transcripts with decreasing abundance in *M. truncatula* and

432    increasing abundance in *C. australe* reinforce the notion that towards the end of seed

433    development, *C. australe* is metabolically active while *M. truncatula* is entering a phase of

434    low metabolic activity and quiescence. Examples of these transcripts are beta-galactosidase

435    (Medtr8g039160), xyloglucan galactosyltransferase (Medtr1g069460) and TCP family TF

436    (Medtr6g015350). Interestingly, these genes lost synteny in *C. australe* compared to their *M.*

437    *truncatula* orthologs. Their involvement in carbohydrate metabolism and control of cell

438    proliferation hints at implication in the germination program. The germination program

439    remains active in *C. australe*, as DS seeds generally do not display developmental arrest,

440    allowing the maintenance of high metabolic activity. In contrast, transcripts of indole-3-acetic

441    acid-amido synthetases, involved in auxin homeostasis, accumulated in *M. truncatula* during

442    development but decreased in *C. australe*. Auxin has been reported to maintain seed

443    dormancy by interacting with ABA (Liu *et al.*, 2013).

444       At the genome level, very few protein-coding genes ($\leq 1\%$) were present in DS

445    species only (Figure 2), supporting the hypothesis that independent evolutionary events gave

446    rise to DS-seeded species. No protein-coding genes were retained in all DT species and lost in

447    all DS species. However, several were lost in all DS and retained in at least half of the DT

448    species. Among these were the transcription factors (TFs) *VERDANDI* and *MYB44-like*.

449    *VERDANDI* participates in ovule identity complex and, when mutated, affects embryo sac

14

Comment citer ce document :
Marques, Costa, M.-C. D., Farrant, J. M., Hilhorst, Buitink, J., Ligterink, Leprince, O.,
Pelletier, S., Gabaldon, T., Schranz, M. E., Delahaie, J., Julca, I., Marcet-Houben, Nijveen, H.,
Derks, Schijlen, E., Zhao, Jonkheer, Chaturi (2019). A blueprint of seed desiccation sensitivity
in the genome of Castanospermum australe. BioRxiv, preprint, . DOI : 10.1101/665661

450    differentiation in *A. thaliana* (Matias-Hernandez *et al.*, 2010; Mendes *et al.*, 2016).

451    Interestingly, *VERDANDI* and *MYB44-like* were also retained in intermediate-seeded species.

452    One gene (*PLAC8*) was lost without retention of paralogs in three out of four DS-

453    seeded species, namely *C. australe*, *Castanea mollissima* and *Elaeis guineensis*. The knock-

454    out of this gene caused increased seed and fruit size in maize (Libault & Stacey, 2010). In

455    addition, the fw2.2 locus containing the *PLAC8* gene has been suggested to be the key to the

456    evolution of tomato fruit size (Frary, 2000). Large seeds and fruits are common features of DS

457    species and presumably reduce the rate of seed drying and hence the risk of desiccation-

458    induced embryo mortality (Daws *et al.*, 2006).

459    Amongst the genes that lost synteny in *C. australe* without detected paralogs, 11 were

460    shared with other DS-seeded species. Two of these genes, *LEA2* and *FIBRILLIN5* accumulate

461    transcripts in *M. truncatula* during seed maturation and upon re-induction of desiccation

462    tolerance in germinated seeds (Terrasson *et al.*, 2013). The gene *GUN5*, a magnesium

463    chelatase involved in retrograde signalling and ABA signalling to the nucleus (Jiang *et al.*,

464    2014), was lost in *C. australe* without paralogs. This pathway is affected in *M. truncatula*

465    *abi5* mutants that produce seeds with strongly reduced longevity (Zinsmeister *et al.*, 2016)

466    and cannot reacquire desiccation tolerance after germination (Terrasson *et al.*, 2013).

467    Examples of genes which lost their ortholog but kept a paralog in *C. australe* are

468    *RETARDED ROOT GROWTH-LIKE* (*RRL)* and *MOTHER OF FT* (*MFT*), involved in ABA-

469    and BR- signalling. *RRL* mediates ABA signal transduction through *ABI4* (Park *et al.*, 2015)

470    and *MFT* regulates seed germination and fertility involving ABA- and BR-signalling

471    pathways (Sun *et al.*, 2010).

472    ABA is involved in the formation of mature DT seeds, and inhibition of their

473    subsequent germination under conditions unfavourable for seedling growth (Finkelstein,

474    2013). During seed development, an increase in ABA content has been related to a transition

475    from growth by cell division to growth by cell enlargement and to cell cycle arrest at the G1/S

476    transition. This increase may be related to the role of ABA in promoting senescence, a process

477    which precedes abscission (Finkelstein, 2013). Additionally, the higher ABA content in the

478    seed coat could play a role in delaying germination of the embryo until ideal conditions for

479    germination are met, or to aid temporal and/or spatial dispersal without immediate loss of

480    viability. Interestingly, the seed coat starts to peel away from the embryo in mature seeds

481    (Figure 1) and this increases markedly once seeds are shed, suggesting increasing lack of

482    inhibition of germination of the embryo.

ABA is also a key regulator of abiotic stress responses and acquisition of desiccation tolerance during seed development (reviewed by (Dekkers *et al.*, 2015)). Disruption of ABA biosynthesis or -signalling leading to lack of or insensitivity to ABA results in loss of seed desiccation tolerance (Verdier *et al.*, 2013). We observed alterations in genes related to ABA signalling, such as *MFT*, *GPCR-TYPE G PROTEIN 1* (*GTG1*), *STH2*, *ABI3* and *ABI5*. *C. australe* lost an ortholog of *MFT* but maintained a paralog. Mutations in this gene may cause ABA hypersensitivity at germination and is associated with dormancy (Vaistij *et al.*, 2013). *STH2* is involved in ABA signalling, is highly expressed during embryogenesis (Xu *et al.*, 2014) and has a high dN/dS in *C. australe*. *GTG1* was lost in *C. mollissima* and its knock-out causes ABA hyposensitivity in *A. thaliana* seeds (Pandey *et al.*, 2009). *ABI3* and *ABI5* showed contrasting expression patterns during *C. australe* seed development compared to *M. truncatula*. These two TFs have been shown to play essential roles in seed development and acquisition of desiccation tolerance (Terrasson *et al.*, 2013; Dekkers *et al.*, 2015; Zinsmeister *et al.*, 2016).

BRs have been implicated in seed development and are known to antagonize seed dormancy and stimulate germination (Steber, 2001). In *C. australe,* several genes involved in BR-biosynthesis and -signalling and seed development have undergone genetic changes (Figure 5). For example, *BRASSINOSTEROID INSENSITIVE-LIKE 3* (*BRL3*) lost an ortholog but kept a paralog; *IKU1* and *MINI3* lost synteny; and *ROT3, DE-ETIOLATED 2* (*DET2*)*, ATBS1* and *BAK1* have higher dN/dS in *C. australe* than in *M. truncatula.* Furthermore, an ortholog of *ATBS1* was highly expressed during late development of *C. australe* contrasting with the decreasing expression in *M. truncatula* seeds. These data support the hypothesis that subcellular metabolism associated with germination is initiated during the late stages of development in the DS seeds of *C. australe*. Overall, our results support the hypothesis that the evolution of desiccation sensitivity was not caused by massive alterations in enzymes and structural proteins but instead by discrete mutations in regulatory genes.

Natural populations often undergo the weakening or removal of a selective force that had been important in the maintenance of a trait, characterizing a "relaxed selection" (Lahti *et al.*, 2009). When a DT-seeded species is subjected to an environment where desiccation tolerance is not an adaptive trait, there should be relaxation of its evolutionary constraints that can eventually lead to its loss. DS-seeded species evolved in environments where the conditions favour immediate germination and seeds are programmed to initiate germination upon, or shortly after shedding (Farrant *et al.*, 1993a; Daws *et al.*, 2006; Berjak & Pammenter, 2013). DT seeds, which very often display a form of dormancy, normally form seed banks in

16

Comment citer ce document :
Marques, Costa, M.-C. D., Farrant, J. M., Hilhorst, Buitink, J., Ligterink, Leprince, O., Pelletier, S., Gabaldon, T., Schranz, M. E., Delahaie, J., Julca, I., Marcet-Houben, Nijveen, H., Derks, Schijlen, E., Zhao, Jonkheer, Chaturi (2019). A blueprint of seed desiccation sensitivity in the genome of Castanospermum australe. BioRxiv, preprint. , DOI : 10.1101/665661

517 the soil. In contrast, DS seeds germinate immediately and usually form seedling banks under

518 shaded forest canopy and take advantage of an eventual light gap for faster establishment.

519 Furthermore, in these species, the generally increased seed size favours seedling

520 establishment under shaded forest conditions (Daws *et al.*, 2006).

521 In summary, seed desiccation sensitivity evolved multiple independent times in

522 environments where water is highly abundant and predictable across long periods, favouring

523 immediate seed germination. In such environments, the evolutionary pressure for DT seeds is

524 relaxed and the production of DS seeds is not disadvantageous. This was the case for *C.*

525 *australe*. Among the currently known DS-non-viviparous seeds, *C. australe* is one of the most

526 sensitive to water loss. We have pinpointed some of the factors behind this sensitivity, namely

527 displacements, loss of synteny and mis-expression of specific genes related to the BR- and

528 ABA-signalling pathways, carbon metabolism, control of cell proliferation, protection against

529 abiotic stresses and modulation of plant stress responses. These alterations are likely to have

530 led to an increased seed size; high starch and low protein and lipid seed content; low

531 accumulation of LEA proteins and RFOs in the seeds; and failure to start the seed maturation

532 drying phase.

533 The low similarity between DS-seeded species confirms the hypothesis that

534 desiccation sensitivity evolved independently. Moreover, it supports the idea that although the

535 evolution of many factors was necessary for the appearance of seed desiccation tolerance,

536 only a few changes in some of these factors are enough for its loss.

537

538 **Author contributions**

539 A.M. and M.-C.D.C. wrote the article; U.C. performed physiological experiments; M.-C.D.C.,

540 E.J., T.Z., M.D. and H.N. performed the bioinformatics; A.M., E.S., M.M.-H., I.J. and T.G.

541 contributed to the genome and transcriptome analysis; J.D. and J.B. performed and analysed

542 the transcriptomics; M.E.S. performed the PacBio sequencing and initial genome analysis; J.B.,

543 H.W.M.H. and J.M.F. initiated and coordinated the work and directed preparation of the article.

544

545 **Acknowledgements**

550 des Pays de la Loire, France (QUALISEM 2009-2013) and the bilateral Partenariat Hubert

551 Curien (PHC) program France–South Africa (grant no. 25903RE) to O.L. and J.B.). We

552 acknowledge David Lalanne and the ANAN platform of the SFR Quasav, Angers, France for

553 the assistance with the microarray analysis. We acknowledge Bas te Lintel Hekkert for library

554 preparation for genome sequencing.

555

556 **References**

557 **Berjak P, Farrant JM, Mycock DJ, Pammenter NW**. **1990**. Recalcitrnat (homoiohydrous)

558 seeds: the enigma of their desiccation-sensitivity. *Seed Science and Technology* **18**: 297–

559 310.

560 **Berjak P, Pammenter NW**. **2000**. What ultrastructure has told us about recalcitrant seeds.

561 *Revista Brasileira de Fisiologia Vegetal* **12**: 22–55.

562 **Berjak P, Pammenter NW**. **2008**. From Avicennia to Zizania: Seed recalcitrance in

563 perspective. *Annals of Botany* **101**: 213–228.

564 **Berjak P, Pammenter NW**. **2013**. Implications of the lack of desiccation tolerance in

565 recalcitrant seeds. *Frontiers in Plant Science* **4**: 1–9.

566 **Bernatzky R, Tanksley SD**. **1986**. Toward a saturated linkage map in tomato based on

567 isozymes and random cDNA sequences. *Genetics* **112**: 887–98.

568 **Bewley JD, Bradford KJ, Hilhorst HWM, Nonogaki H**. **2013**. *Seeds. Physiology of*

569 *development, germination and dormancy* (JD Bewley, KJ Bradford, HWM Hilhorst, and H

570 Nonogaki, Eds.). Springer.

571 **Boetzer M, Pirovano W**. **2014**. SSPACE-LongRead: scaffolding bacterial draft genomes

572 using long read sequence information. *BMC Bioinformatics* **15**: 211.

573 **Cardoso D, de Queiroz LP, Pennington RT, de Lima HC, Fonty É, Wojciechowski MF,**

574 **Lavin M**. **2012**. Revisiting the phylogeny of papilionoid legumes: New insights from

575 comprehensively sampled early-branching lineages. *American Journal of Botany* **99**:

576 1991–2013.

577 **Chaisson MJ, Tesler G**. **2012**. Mapping single molecule sequencing reads using basic local

578 alignment with successive refinement (BLASR): application and theory. *BMC*

579 *Bioinformatics* **13**: 238.

580 **Chatelain E, Hundertmark M, Leprince O, Gall S Le, Satour P, Deligny-Penninck S,**

581 **Rogniaux H, Buitink J**. **2012**. Temporal profiling of the heat-stable proteome during late

582 maturation of *Medicago truncatula* seeds identifies a restricted subset of late

18

583     embryogenesis abundant proteins associated with longevity. *Plant, Cell and Environment*

584     **35**: 1440–1455.

585 **Costa M-CD, Artur MAS, Maia J, Jonkheer E, Derks MFL, Nijveen H, Williams B,**

586     **Mundree SG, Jiménez-Gómez JM, Hesselink T, *et al.* 2017**. A footprint of desiccation

587     tolerance in the genome of *Xerophyta viscosa*. *Nature Plants* **3**: 17038.

588 **Dang NX, Popova A V., Hundertmark M, Hincha DK. 2014**. Functional characterization

589     of selected LEA proteins from *Arabidopsis thaliana* in yeast and in vitro. *Planta* **240**: 325–

590     336.

591 **Daws MI, Garwood NC, Pritchard HW. 2006**. Prediction of desiccation sensitivity in seeds

592     of woody species: A probabilistic model based on two seed traits and 104 species. *Annals*

593     *of Botany* **97**: 667–674.

594 **Dekkers BJW, Costa MCD, Maia J, Bentsink L, Ligterink W, Hilhorst HWM. 2015**.

595     Acquisition and loss of desiccation tolerance in seeds: from experimental model to

596     biological relevance. *Planta* **241**: 563–577.

597 **Delahaie J, Hundertmark M, Bove J, Leprince O, Rogniaux H, Buitink J. 2013**. LEA

598     polypeptide profiling of recalcitrant and orthodox legume seeds reveals ABI3-regulated

599     LEA protein abundance linked to desiccation tolerance. *Journal of Experimental Botany*

600     **64**: 4559–4573.

601 **Djemel N, Guedon D, Lechevalier A, Salon C, Miquel M, Prosperi JM, Rochat C, Boutin**

602     **JP. 2005**. Development and composition of the seeds of nine genotypes of the *Medicago*

603     *truncatula* species complex. *Plant Physiology and Biochemistry* **43**: 557–566.

604 **Eddy SR. 2011**. Accelerated profile HMM searches. *PLoS Computational Biology* **7**.

605 **English AC, Richards S, Han Y, Wang M, Vee V, Qu J, Qin X, Muzny DM, Reid JG,**

606     **Worley KC, *et al.* 2012**. Mind the gap: Upgrading genomes with Pacific Biosciences RS

607     Long-Read sequencing technology. *PLoS ONE* **7**: 1–12.

608 **Farrant JM, Berjak P, Cutting JGM, Pammenter NW. 1993a**. The role of plant growth

609     regulators in the development and germination of the desiccation sensitive recalcitrant

610     seeds of *Avicennia marina*. *Seed Science Research* **3**: 55–63.

611 **Farrant JM, Berjak P, Pammenter NW. 1992**. Proteins in development and germination of

612     a desiccation sensitive (recalcitrant) seed species. *Plant Growth Regulation* **11**: 257–265.

613 **Farrant JM, Cooper K, Nell H. 2012**. Desiccation tolerance. *Plant Stress Physiology*: 238–

614     265.

615 **Farrant JM, Pammenter NW, Berjak P. 1993b**. Seed development in relation to

616     desiccation tolerance: A comparison between desiccation-sensitive (recalcitrant) seeds of

617    *Avicennia marina* and desiccation-tolerant types. *Seed Science Research* **3**.

618  **Finkelstein R**. **2013**. Abscisic acid synthesis and response. *The Arabidopsis Book* **11**: e0166.

619  **Finn RD, Clements J, Eddy SR**. **2011**. HMMER web server: Interactive sequence similarity
620    searching. *Nucleic Acids Research* **39**: 1–9.

621  **Floková K, Tarkowská D, Miersch O, Strnad M, Wasternack C, Novák O**. **2014**.
622    UHPLC-MS/MS based target profiling of stress-induced phytohormones. *Phytochemistry*
623    **105**: 147–157.

624  **Francini A, Galleschi L, Saviozzi F, Pinzino C, Izzo R, Sgherri C, Navari-Izzo F**. **2006**.
625    Enzymatic and non-enzymatic protective mechanisms in recalcitrant seeds of *Araucaria*
626    *bidwillii* subjected to desiccation. *Plant Physiology and Biochemistry* **44**: 556–563.

627  **Frary A**. **2000**. fw2.2: A quantitative trait locus key to the evolution of tomato fruit size.
628    *Science* **289**: 85–88.

629  **Gurevich A, Saveliev V, Vyahhi N, Tesler G**. **2013**. QUAST: quality assessment tool for
630    genome assemblies. *Bioinformatics* **29**: 1072–1075.

631  **Hamilton KN, Offord CA, Cuneo P, Deseo MA**. **2013**. A comparative study of seed
632    morphology in relation to desiccation tolerance and other physiological responses in 71
633    Eastern Australian rainforest species. *Plant Species Biology* **28**: 51–62.

634  **Holt C, Yandell M**. **2011**. MAKER2: an annotation pipeline and genome-database
635    management tool for second-generation genome projects. *BMC Bioinformatics* **12**: 491.

636  **Huerta-Cepas J, Capella-Gutiérrez S, Pryszcz LP, Marcet-Houben M, Gabaldón T**.
637    **2014**. PhylomeDB v4: zooming into the plurality of evolutionary histories of a genome.
638    *Nucleic Acids Research* **42**: D897–D902.

639  **Huerta-Cepas J, Gabaldon T**. **2011**. Assigning duplication events to relative temporal scales
640    in genome-wide studies. *Bioinformatics* **27**: 38–45.

641  **Hundertmark M, Hincha DK**. **2008**. LEA (Late Embryogenesis Abundant) proteins and
642    their encoding genes in *Arabidopsis thaliana*. *BMC Genomics* **9**: 118.

643  **Jiang S-C, Mei C, Wang X-F, Zhang D-P**. **2014**. A hub for ABA signaling to the nucleus:
644    Significance of a cytosolic and nuclear dual-localized PPR protein SOAR1 acting
645    downstream of Mg-chelatase H subunit. *Plant Signaling & Behavior* **9**: e972899.

646  **Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, Salzberg SL**.
647    **2004**. Versatile and open software for comparing large genomes. *Genome Biology* **5**: R12.

648  **Lahti DC, Johnson NA, Ajie BC, Otto SP, Hendry AP, Blumstein DT, Coss RG,
649    Donohue K, Foster SA**. **2009**. Relaxed selection in the wild. *Trends in Ecology &*
650    *Evolution* **24**: 487–496.

651   **Langmead B, Salzberg SL**. **2012**. Fast gapped-read alignment with Bowtie 2. *Nature*
652       *methods* **9**: 357–9.

653   **Leprince O, Pellizzaro A, Berriri S, Buitink J**. **2017**. Late seed maturation: Drying without
654       dying. *Journal of Experimental Botany* **68**: 827–841.

655   **Libault M, Stacey G**. **2010**. Evolution of FW2.2-like (FWL) and PLAC8 genes in
656       eukaryotes. *Plant Signaling & Behavior* **5**: 1226–1228.

657   **Liu S, Wang X, Wang H, Xin H, Yang X, Yan J, Li J, Tran LSP, Shinozaki K,**
658       **Yamaguchi-Shinozaki K, *et al.* 2013**. Genome-wide analysis of ZmDREB genes and their
659       association with natural variation in drought tolerance at seedling stage of *Zea mays* L.
660       *PLoS Genetics* **9**.

661   **Luo M, Dennis ES, Berger F, Peacock WJ, Chaudhury A**. **2005**. MINISEED3 (MINI3), a
662       WRKY family gene, and HAIKU2 (IKU2), a leucine-rich repeat (LRR) KINASE gene, are
663       regulators of seed size in Arabidopsis. *Proceedings of the National Academy of Sciences*
664       **102**: 17531–17536.

665   **Maere S, Heymans K, Kuiper M**. **2005**. BiNGO: A Cytoscape plugin to assess
666       overrepresentation of Gene Ontology categories in Biological Networks. *Bioinformatics*
667       **21**: 3448–3449.

668   **Marques A, Buijs G, Ligterink W, Hilhorst H**. **2018**. Evolutionary ecophysiology of seed
669       desiccation sensitivity. *Functional Plant Biology* **45**: 1083.

670   **Matias-Hernandez L, Battaglia R, Galbiati F, Rubes M, Eichenberger C, Grossniklaus**
671       **U, Kater MM, Colombo L**. **2010**. VERDANDI is a direct target of the MADS domain
672       ovule identity complex and affects embryo sac differentiation in arabidopsis. *The Plant*
673       *Cell* **22**: 1702–1715.

674   **Mendes MA, Guerra RF, Castelnovo B, Silva-Velazquez Y, Morandini P, Manrique S,**
675       **Baumann N, Groß-Hardt R, Dickinson H, Colombo L**. **2016**. Live and let die: a REM
676       complex promotes fertilization through synergid cell death in Arabidopsis. *Development*
677       **143**: 2780–2790.

678   **Mugal CF, Wolf JBW, Kaj I**. **2014**. Why time matters: Codon evolution and the temporal
679       dynamics of dN/dS. *Molecular Biology and Evolution* **31**: 212–231.

680   **Pandey S, Nelson DC, Assmann SM**. **2009**. Two novel GPCR-Type G proteins are abscisic
681       acid receptors in arabidopsis. *Cell* **136**: 136–148.

682   **Park S-Y, Peterson FC, Mosquna A, Yao J, Volkman BF, Cutler SR**. **2015**. Agrochemical
683       control of plant water use using engineered abscisic acid receptors. *Nature* **520**: 545–548.

684   **Righetti K, Vu JL, Pelletier S, Vu BL, Glaab E, Lalanne D, Pasha A, Patel R V., Provart**

21

685     **NJ, Verdier J,** *et al.* **2015**. Inference of longevity-related genes from a robust

686     coexpression network of seed maturation identifies regulators linking seed storability to

687     biotic defense-related pathways. *The Plant Cell* **27**: tpc.15.00632.

688     **Rosnoblet C, Aubry C, Leprince O, Vu BL, Rogniaux H, Buitink J**. **2007**. The regulatory

689     gamma subunit SNF4b of the sucrose non-fermenting-related kinase complex is involved

690     in longevity and stachyose accumulation during maturation of *Medicago truncatula* seeds.

691     *Plant Journal* **51**: 47–59.

692     **Rosvall M, Bergstrom CT**. **2008**. Maps of random walks on complex networks reveal

693     community structure. *Proceedings of the National Academy of Sciences* **105**: 1118–1123.

694     **Schranz ME, Mohammadin S, Edger PP**. **2012**. Ancient whole genome duplications,

695     novelty and diversification: the WGD Radiation Lag-Time Model. *Current Opinion in*

696     *Plant Biology* **15**: 147–153.

697     **Sershen, Berjak P, Pammenter NW, Wesley-Smith J**. **2012**. The effects of various

698     parameters during processing for cryopreservation on the ultrastructure and viability of

699     recalcitrant zygotic embryos of *Amaryllis belladonna*. *Protoplasma* **249**: 155–169.

700     **Simão FA, Waterhouse RM, Ioannidis P, Kriventseva E V., Zdobnov EM**. **2015**.

701     BUSCO: assessing genome assembly and annotation completeness with single-copy

702     orthologs. *Bioinformatics* **31**: 3210–3212.

703     **Song L, Florea L, Langmead B**. **2014**. Lighter: fast and memory-efficient sequencing error

704     correction without counting. *Genome Biology* **15**: 509.

705     **Steadman KJ, Pritchard HW, Dey PM**. **1996**. Tissue-specific soluble sugars in seeds as

706     indicators of storage category. *Annals of Botany* **77**: 667–674.

707     **Steber CM**. **2001**. A role for brassinosteroids in germination in arabidopsis. *Plant Physiology*

708     **125**: 763–769.

709     **Sun Y, Fan XY, Cao DM, Tang W, He K, Zhu JY, He JX, Bai MY, Zhu S, Oh E,** *et al.*

710     **2010**. Integration of brassinosteroid signal transduction with the transcription network for

711     plant growth regulation in arabidopsis. *Developmental Cell* **19**: 765–777.

712     **Terrasson E, Buitink J, Righetti K, Ly Vu B, Pelletier S, Zinsmeister J, Lalanne D,**

713     **Leprince O**. **2013**. An emerging picture of the seed desiccome: confirmed regulators and

714     newcomers identified using transcriptome comparison. *Frontiers in Plant Science* **4**: 1–16.

715     **Tunnacliffe A, Wise MJ**. **2007**. The continuing conundrum of the LEA proteins.

716     *Naturwissenschaften* **94**: 791–812.

717     **Vaistij FE, Gan Y, Penfield S, Gilday AD, Dave A, He Z, Josse E-M, Choi G, Halliday**

718     **KJ, Graham IA**. **2013**. Differential control of seed primary dormancy in Arabidopsis

22

719 ecotypes by the transcription factor SPATULA. *Proceedings of the National Academy of*
720 *Sciences* **110**: 10866–10871.

721 **VanBuren R, Wai CM, Zhang Q, Song X, Edger PP, Bryant D, Michael TP, Mockler**
722 **TC, Bartels D**. **2017**. Seed desiccation mechanisms co-opted for vegetative desiccation in
723 the resurrection grass *Oropetium thomaeum*. *Plant Cell and Environment* **40**: 2292–2306.

724 **Verdier J, Lalanne D, Pelletier S, Torres-Jerez I, Righetti K, Bandyopadhyay K,**
725 **Leprince O, Chatelain E, Vu BL, Gouzy J,** *et al.* **2013**. A regulatory network-based
726 approach dissects late maturation processes related to the acquisition of desiccation
727 tolerance and longevity of *Medicago truncatula* seeds. *Plant Physiology* **163**: 757–774.

728 **Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng**
729 **Q, Wortman J, Young SK,** *et al.* **2014**. Pilon: An integrated tool for comprehensive
730 microbial variant detection and genome assembly improvement. *PLoS ONE* **9**: e112963.

731 **Wang Y, Tang H, DeBarry JD, Tan X, Li J, Wang X, Lee T -h., Jin H, Marler B, Guo H,**
732 *et al.* **2012**. MCScanX: a toolkit for detection and evolutionary analysis of gene synteny
733 and collinearity. *Nucleic Acids Research* **40**: e49–e49.

734 **Wang M, Verdier J, Benedito VA, Tang Y, Murray JD, Ge Y, Becker JD, Carvalho H,**
735 **Rogers C, Udvardi M,** *et al.* **2013**. LegumeGRN: A gene regulatory network prediction
736 server for functional and comparative studies. *PLoS ONE* **8**.

737 **Xu D, Li J, Gangappa SN, Hettiarachchi C, Lin F, Andersson MX, Jiang Y, Deng XW,**
738 **Holm M**. **2014**. Convergence of light and ABA signaling on the ABI5 promoter (L-J Qu,
739 Ed.). *PLoS Genetics* **10**: e1004197.

740 **Yang Z**. **2007**. PAML 4: Phylogenetic analysis by maximum likelihood. *Molecular Biology*
741 *and Evolution* **24**: 1586–1591.

742 **Ye C, Hill CM, Wu S, Ruan J, Ma Z**. **2016**. DBG2OLC: Efficient assembly of large
743 genomes using long erroneous reads of the third generation sequencing technologies.
744 *Scientific Reports* **6**: 31900.

745 **Ye C, Ma ZS, Cannon CH, Pop M, Yu DW**. **2012**. Exploiting sparseness in de novo
746 genome assembly. *BMC Bioinformatics* **13**: S1.

747 **Yin Y, Cheong H, Friedrichsen D, Zhao Y, Hu J, Mora-Garcia S, Chory J**. **2002**. A
748 crucial role for the putative Arabidopsis topoisomerase VI in plant growth and
749 development. *Proceedings of the National Academy of Sciences* **99**: 10191–10196.

750 **Zhao T, Holmer R, Bruijn S de, Angenent GC, van den Burg HA, Schranz ME**. **2017**.
751 Phylogenomic synteny network analysis of MADS-box transcription factor genes reveals
752 lineage-specific transpositions, ancient tandem duplications, and deep positional

Version preprint

753    conservation. *The Plant Cell*: tpc.00312.2017.

754    **Zinsmeister J, Lalanne D, Terrasson E, Chatelain E, Vandecasteele C, Vu BL, Dubois-**

755    **Laurent C, Geoffriau E, Signor C Le, Dalmais M, *et al.* 2016**. ABI5 is a regulator of

756    seed maturation and longevity in legumes. *The Plant Cell* **28**: 2735–2754.

757

758    **Supporting Information**

759

760    **Figure S1.** Hierarchical clustering of expression values.

761

762    **Figure S2.** Phylogenetic tree showing duplication rates of selected species.

763

764    **Table S1**: Gene ontology (GO) enrichment analysis of biological processes in relation to the

765    acquisition of tolerance to water loss and to seed maturation in *Medicago truncatula* and

766    *Castanospermum australe*.

767

768    **Table S2.** Genes changing transcript abundance in *Castanospermum australe* and *Medicago*

769    *truncatula* in comparable seed developmental stages.

770

771    **Table S3.** Genes changing transcript abundance in *Castanospermum australe* and *Medicago*

772    *truncatula* in comparable seed developmental stages during final maturation.

773

774    **Table S4.** Protein-coding genes lost in all DS and retained in at least half of the DT species.

775

776    **Table S5.** *Castanospermum australe* protein-coding genes with dN/dS (number of

777    nonsynonymous substitutions per non-synonymous site (dN) in a given period of time divided

778    by the number of synonymous substitutions per synonymous site (dS) in the same period)

779    ratio ≥2.

780

781    **Figure legends**

782

783    **Figure 1. *Castanospermum australe* seed and pod developmental stages.** MAF: months

784    after flowering.

785 **Figure 2. Phylome.** Reconstruction of the phylome was based on a concatenated alignment of

786 183 single copy proteins that are present in at least 19 out of the 20 species surveyed. Species

787 names are coloured according to their seed storage category.

788 **Figure 3**. **Synonymous mutations of duplicated genes and syntenic blocks.** Histograms of

789 synonymous mutations (Ks) of duplicated genes (A-F) and average Ks of syntenic blocks (G-

790 L). (A-B) *Castanospermum australe*, (C-D) *Lotus japonicus*, (E-F) *Medicago truncatula*, (G-

791 H) *Glycine max*, (I-J) *Phaseolus vulgaris*, (K-L) *Trifolium pratense.*

792 **Figure 4. Hypothetical model for brassinosteroid-regulated seed development (adapted**

793 **from Jiang et al. (2013)).** Red shapes indicate genes that lost synteny in *Castanospermum*

794 *australe* compared to *Medicago truncatula*. Green shapes indicate genes with higher dN/dS in

795 *C. australe* than in *M. truncatula*. Blue shapes indicate genes that lost an ortholog, but kept a

796 paralog in *C. australe*. BR: brassinosteroid. ATBSI1: *ACTIVATION-TAGGED BRI1*

797 *(BRASSINOSTEROID-INSENSITIVE1)-SUPPRESSOR1*. AP2: *APETALA 2*. ARF2: *AUXIN*

798 *RESPONSE FACTOR 2*. BAK1: *BRI1-ASSOCIATED RECEPTOR KINASE 1*. BRI1:

799 *BRASSINOSTEROID INSENSITIVE 1*. BRL3: *BRASSINOSTEROID INSENSITIVE-LIKE 3*.

800 BZR1: *BRASSINAZOLE RESISTANT 1*. DET2: *DE-ETIOLATED 2*. IKU: *HAIKU*. MINI3:

801 *MINISEED 3*. ROT3: *ROTUNDIFOLIA 3*. SHB1: *SHOEBOX 1*.

802

803 **Table legends**

804 **Table 1.** Phenotypic parameters associated with late seed development of *Castanospermum*

805 *australe.* Point of attachment refers to tissue attaching the embryo to the seed pod.

806 **Table 2.** Overview of assembly, annotation, polymorphism and repeat elements on the

807 *Castanospermum australe* genome. N50: scaffold size above which 50% of the total length of

808 the sequence assembly can be found. L90: number of contigs whose summed length contains

809 at least 90% of the sum of the total length of the sequence assembly. rRNA: ribosomal RNA.

810 snRNA: small nuclear RNA. tRNA: transfer RNA. SNP: single nucleotide polymorphism.

811 INDEL: insertion or deletion of bases in the DNA. SINE: short interspersed nuclear element.

812 LINE: long interspersed nuclear element. LTR: long terminal repeat. L1: LINE-1. RTE:

813 retrotransposable element. *hAT*: *hobo/Ac/Tam3*. PIF: P Instability Factor. CMC:

814 CACTA/Mirage/Chapaev.

815

816

Comment citer ce document :
Marques, Costa, M.-C. D., Farrant, J. M., Hilhorst, Buitink, J., Ligterink, Leprince, O.,
Pelletier, S., Gabaldon, T., Schranz, M. E., Delahaie, J., Julca, I., Marcet-Houben, Nijveen, H.,
Derks, Schijlen, E., Zhao, Jonkheer, Chaturi (2019). A blueprint of seed desiccation sensitivity
in the genome of Castanospermum australe. BioRxiv, preprint. . DOI : 10.1101/665661
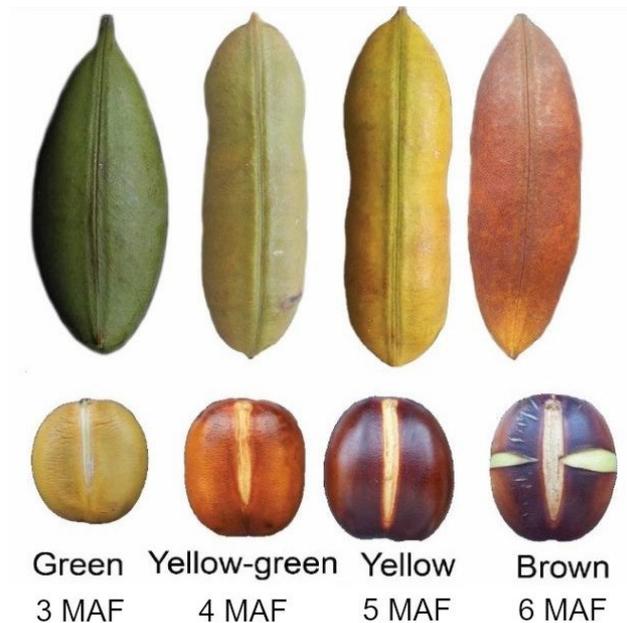
**Figures**
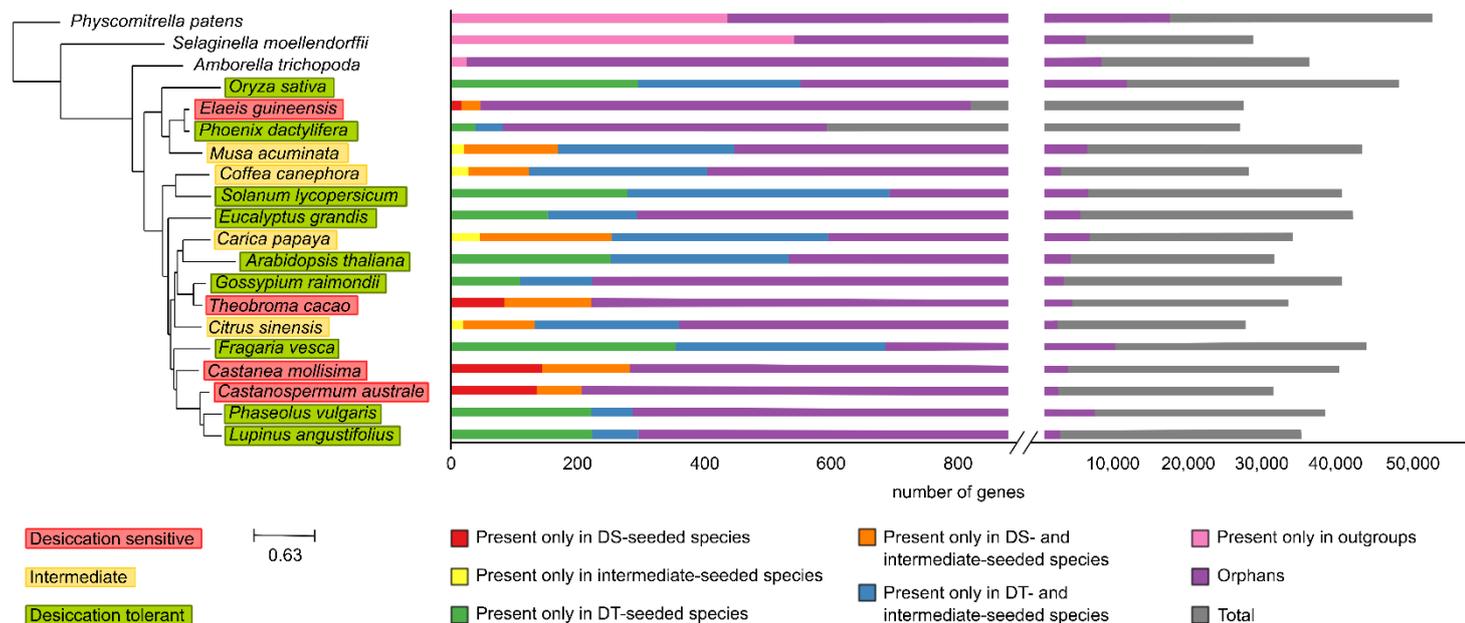
818

819

820

821

822

823

824

825

826

827

828

829



830 **Figure 1.** *Castanospermum australe* **seed and pod developmental stages.** MAF: months
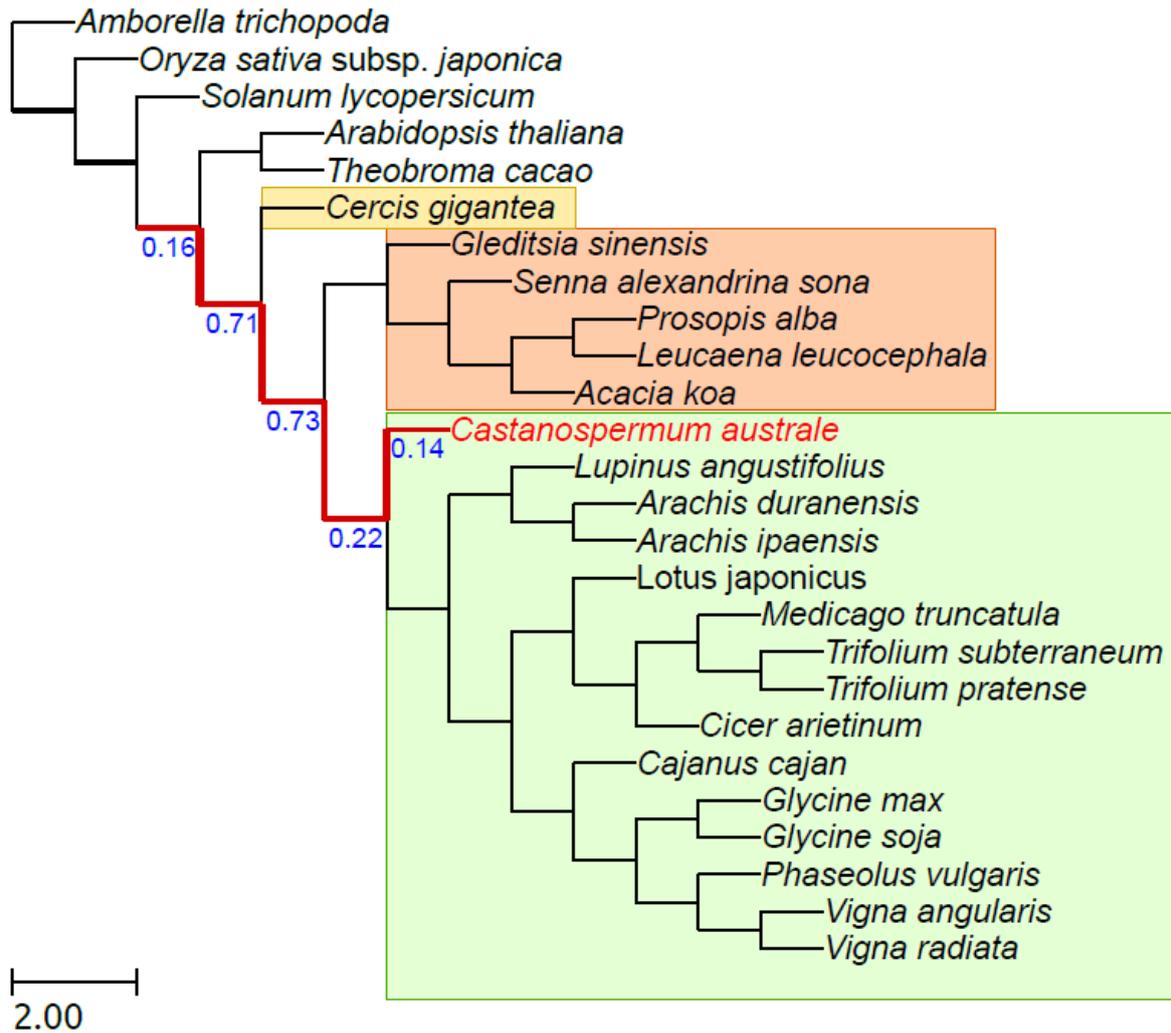
831 after flowering.

832
833



845 **Figure 2. Phylome.** Reconstruction of the phylome was based on a concatenated alignment of

846 183 single copy proteins that are present in at least 19 out of the 20 species surveyed. Species
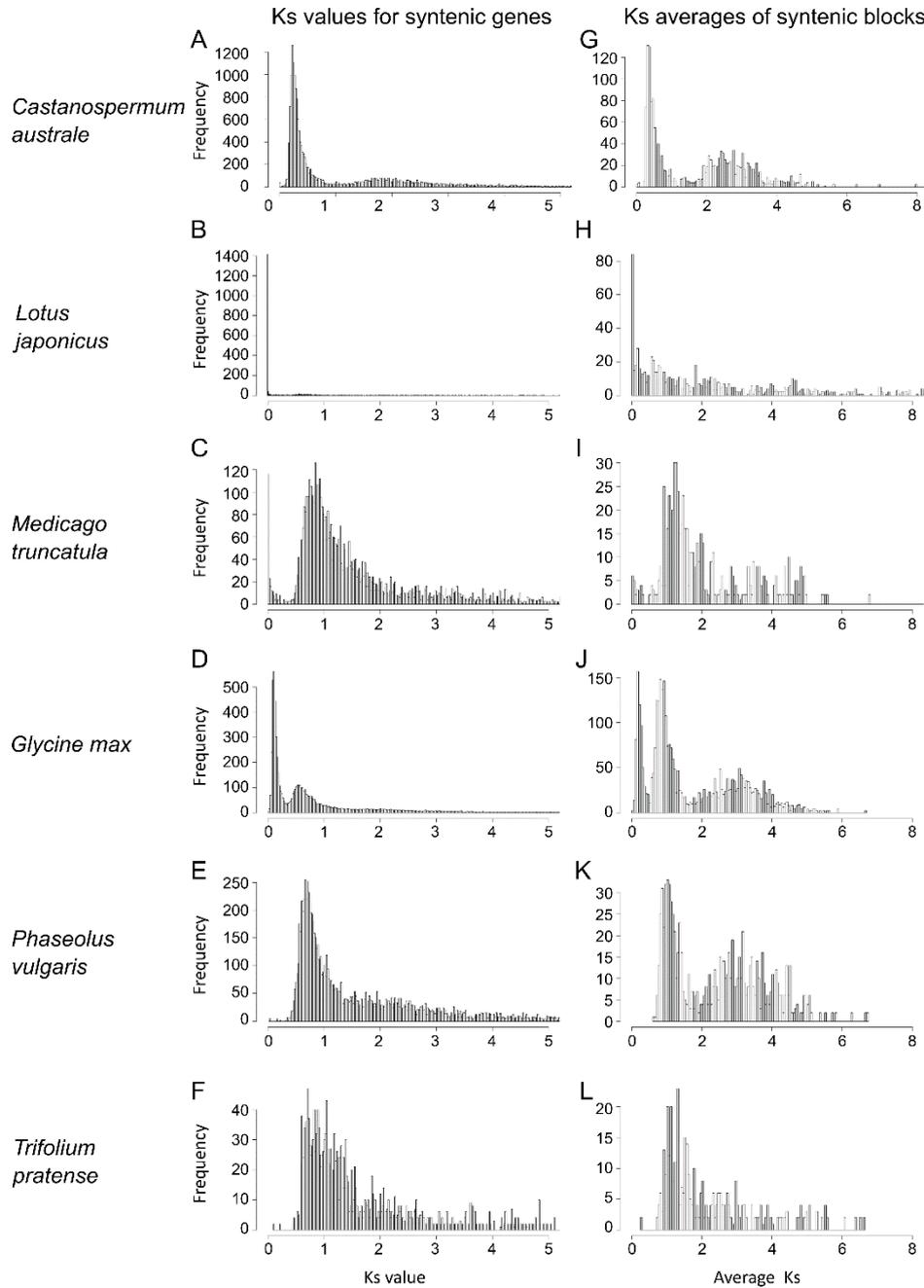
847 names are coloured according to their seed storage category.

**Figure 3.** Phylogenetic tree showing duplication rates of selected species including species of the subfamilies: Mimosoideae (inside yellow rectangle) and Caesalpinioideae (inside orange rectangle) and Papilionoideae (inside green rectangle). Phylogentic tree based on: Lavin, M., Herendeen, P.S. and Wojciechowski, M.F. (2005) evolutionary rates analysis of leguminosae implicates a rapid diversification of lineages during the tertiary. Systematic Biology 54, 575–594.
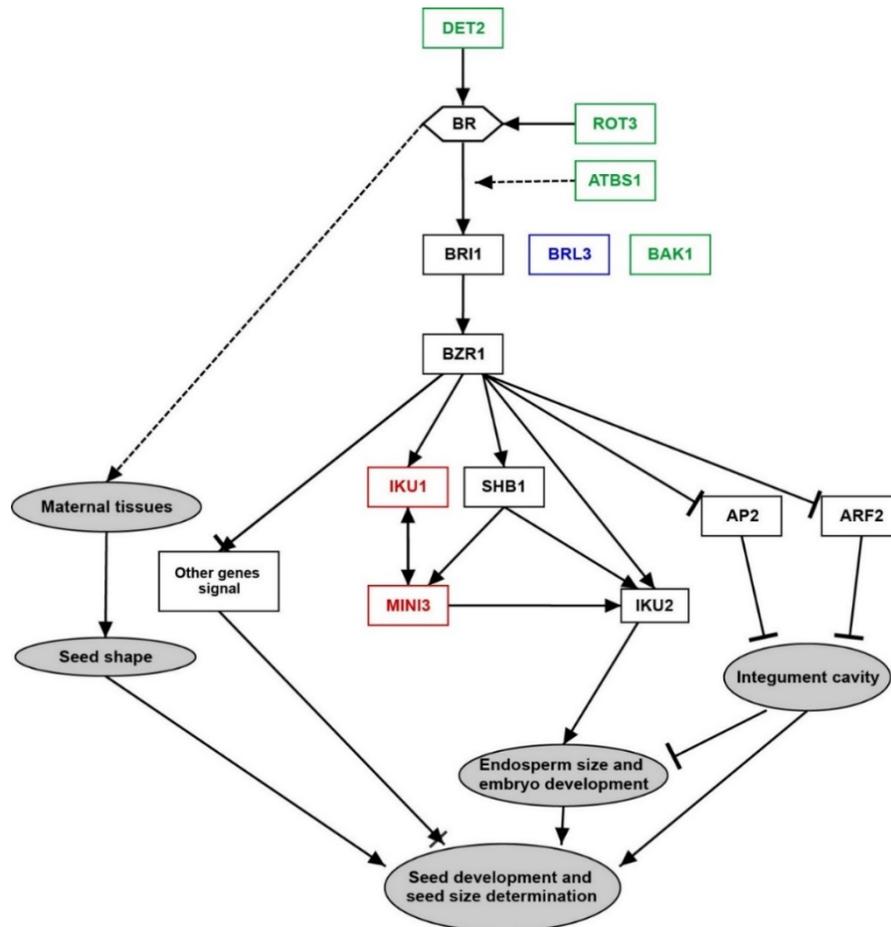
856

857



858

**Figure 4**. **Synonymous mutations of duplicated genes and syntenic blocks.** Histograms of
synonymous mutations (Ks) of duplicated genes (A-F) and average Ks of syntenic blocks (G-
L). (A-B) *Castanospermum australe*, (C-D) *Lotus japonicus*, (E-F) *Medicago truncatula*, (G-
H) *Glycine max*, (I-J) *Phaseolus vulgaris*, (K-L) *Trifolium pratense.*

863

29

865

**Figure 5. Hypothetical model for brassinosteroid-regulated seed development (adapted from Jiang et al. (2013)).** Red shapes indicate genes that lost synteny in *Castanospermum australe* compared to *Medicago truncatula*. Green shapes indicate genes with higher dN/dS in *C. australe* than in *M. truncatula*. Blue shapes indicate genes that lost an ortholog, but kept a paralog in *C. australe*. BR: brassinosteroid. ATBSI1: *ACTIVATION-TAGGED BRI1 (BRASSINOSTEROID-INSENSITIVE1)-SUPPRESSOR1*. AP2: *APETALA 2*. ARF2: *AUXIN RESPONSE FACTOR 2*. BAK1: *BRI1-ASSOCIATED RECEPTOR KINASE 1*. BRI1: *BRASSINOSTEROID INSENSITIVE 1*. BRL3: *BRASSINOSTEROID INSENSITIVE-LIKE 3*. BZR1: *BRASSINAZOLE RESISTANT 1*. DET2: *DE-ETIOLATED 2*. IKU: *HAIKU*. MINI3: *MINISEED 3*. ROT3: *ROTUNDIFOLIA 3*. SHB1: *SHOEBOX 1*.

30

877 **Table 1.** Phenotypic parameters associated with late seed development of *Castanospermum australe.* Point of attachment refers to tissue attaching the embryo to the seed pod.

| Pod stage | | Green | | Yellow-green | | Yellow | | Brown | |
|---|---|---|---|---|---|---|---|---|---|
| | | average | SE | average | SE | average | SE | average | SE |
| Whole seed mass (g) | | 19.236 | 0.586 | 29.618 | 0.959 | 38.037 | 0.867 | 40.389 | 0.720 |
| Seed coat | mass (g) | 2.496 | 0.082 | 1.342 | 0.039 | 1.210 | 0.032 | 1.228 | 0.037 |
| | water content ($gH_2O$ $g^{-1}$ dry mass) | 2.042 | 0.265 | 3.065 | 0.146 | 1.001 | 0.055 | 0.652 | 0.045 |
| Germination (%) of axes on MS media | | 93 | | 100 | | 100 | | 100 | |
| Time (days) to reach 50% germination | | 25 | | 8 | | 8 | | 8 | |
| Water content ($gH_2O$ $g^{-1}$ dry mass) | axes | 4.030 | 0.149 | 3.458 | 0.113 | 2.412 | 0.061 | 2.426 | 0.045 |
| | cotyledons | 5.084 | 0.209 | 3.255 | 0.231 | 1.559 | 0.073 | 1.783 | 0.122 |
| Water content ($gH_2O$ $g^{-1}$ dry mass) at which 50% loss of viability occurs | axes | 0.448 | 0.021 | 0.342 | 0.016 | 0.228 | 0.012 | 0.221 | 0.009 |
| | cotyledons | 0.820 | 0.026 | 0.735 | 0.021 | 0.766 | 0.021 | 0.700 | 0.027 |
| Glucose (mg $g^{-1}$ dry mass) | axes | 0.391 | 0.144 | 1.545 | 0.194 | 1.162 | 0.366 | 1.088 | 0.410 |
| | cotyledons | 0.829 | 0.719 | 0.528 | 0.228 | 0.261 | 0.106 | 0.661 | 0.064 |
| | point of attachment | | | 0.515 | 1.486 | 0.245 | 0.193 | 0.450 | 0.257 |
| Fructose (mg $g^{-1}$ dry mass) | axes | 0.243 | 0.218 | 0.648 | 0.063 | 0.912 | 0.456 | 0.100 | 0.173 |
| | cotyledons | 0.491 | 0.438 | 0.343 | 0.069 | 0.336 | 0.033 | 0.359 | 0.115 |
| | point of attachment | | | 0.950 | 1.039 | 7.259 | 11.653 | 0.781 | 0.302 |
| Sucrose (mg $g^{-1}$ dry mass) | axes | 40.615 | 6.792 | 81.402 | 3.842 | 89.790 | 13.852 | 82.760 | 6.089 |
| | cotyledons | 41.428 | 18.578 | 25.854 | 6.820 | 63.297 | 3.108 | 58.342 | 12.720 |
| | point of attachment | | | 1.858 | 0.515 | 0.783 | 0.319 | 1.434 | 1.312 |
| Stachyose (mg $g^{-1}$ dry mass) | axes | 1.947 | 0.644 | 1.767 | 0.179 | 1.139 | 0.116 | 1.043 | 0.170 |
| | cotyledons | 2.948 | 1.705 | 0.154 | 0.135 | 0.112 | 0.043 | 0.128 | 0.026 |
| | point of attachment | | | 1.188 | 1.858 | 1.948 | 0.167 | 1.288 | 0.281 |
| Raffinose (mg $g^{-1}$ dry mass) | axes | 1.990 | 0.217 | 0.546 | 0.184 | 0.000 | 0.000 | 0.205 | 0.354 |
| | cotyledons | 1.576 | 0.309 | 0.328 | 0.567 | 0.639 | 0.619 | 0.531 | 0.477 |
| | point of attachment | | | 0.000 | 0.950 | 0.000 | 0.000 | 0.000 | 0.000 |
| ABA content (pmol $g^{-1}$ dry mass) | axes | 8556.94 | 62.50 | 7676.37 | 42.89 | 817.89 | 82.08 | 444.95 | 7.65 |
| | cotyledons | 18393.12 | 335.31 | 16153.81 | 160.86 | 4311.68 | 188.76 | 1860.64 | 35.90 |
| | point of attachment | 9404.77 | 101.66 | 11710.72 | 265.05 | 19400.48 | 181.23 | 24172.84 | 574.80 |
| | seed coat | 9175.61 | 330.12 | 32259.10 | 115.59 | 15779.45 | 125.60 | 23380.20 | 216.08 |

**Table 2.** Overview of assembly, annotation, polymorphism and repeat elements on the *Castanospermum australe* genome. N50: scaffold size above which 50% of the total length of the sequence assembly can be found. L90: number of contigs whose summed length contains at least 90% of the sum of the total length of the sequence assembly. rRNA: ribosomal RNA. snRNA: small nuclear RNA. tRNA: transfer RNA. SNP: single nucleotide polymorphism. INDEL: insertion or deletion of bases in the DNA. SINE: short interspersed nuclear element. LINE: long interspersed nuclear element. LTR: long terminal repeat. L1: LINE-1. RTE: retrotransposable element. *hAT*: *hobo/Ac/Tam3*. PIF: P Instability Factor. CMC: CACTA/Mirage/Chapaev.

| Assembly | Number | N50 (Kb) | L90 (kb) | Total length | Alignment rate |
|---|---|---|---|---|---|
| Contigs | 1,210 | 761.1 | - | - | - |
| Scaffolds | 1,027 | 832.6 | 495 | 381.7 Mb | 97.7% |

| Annotation | Number | Mean lengh (bp) | Density | Genome percentage |
|---|---|---|---|---|
| Protein coding genes | 29,124 | 4814.9 | - | 36.7% |
| Exons | 180,329 | 232.7 | 6.2 exons/gene | 11.1 % |
| Introns | 141,174 | 696.1 | 5.0 introns/gene | 25.6% |
| rRNA | 149 | 623.83 | - | 0.0% |
| snRNA | 60 | 121.72 | - | 0.0% |
| tRNA | 310 | 75.55 | - | 0.0% |
| Transposable elements | 110,949 | - | - | 15.5% |

| Polymorphisms | Number | Density |
|---|---|---|
| SNPs | 352,963 | $0.92$ kb$^{-1}$ |
| INDELS | 328,918 | $0.86$ kb$^{-1}$ |
| Multi-allelic sites | 6,328 | $0.02$ kb$^{-1}$ |

| Repeat genus | Length | Abundant species |
|---|---|---|
| SINE | 118,770 (0.05%) | - |
| LINE | 3,834,244 (1.6%) | L1, RTE |
| DNA transposon | 11,400,282 (4.8%) | *hAT*, PIF, CMC |
| LTR | 42,305,056 (17.84) | Gypsy, Copia |
| Unclassified/Simple | 177,648,468 (74.9%) | - |

878

32

Version preprint

33